

Exceptional service in the national interest

FORESTR

Searching for patterns in random forests

Katherine Goode Joint work with J. Derek Tucker

March 12, 2025



Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525

S A N D 2 0 2 5 - 0 2 8 6 2 P E







Background: Sandia & Machine Learning



FORESTR: Topology Patterns in Random Forests

BACKGROUND

SANDIA & MACHINE LEARNING

SANDIA IS A FEDERALLY FUNDED RESEARCH AND DEVELOPMENT CENTER MANAGED AND OPERATED BY

National Technology & Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc.

Government-owned, contractor-operated

FFRDCs are long-term strategic partners to the federal government, operating in the public interest with objectivity and independence and maintaining core competencies in missions of national significance GJ

WE HAVE FACILITIES **ACROSS THE NATION**



NATIONAL SECURITY IS OUR MISSION

Sandia delivers essential science and technology to address the nation's most challenging security issues





PURPOSE	We render exceptional service in the national interest
VISION	We make Sandia a leader in keeping the world safe and secure
MISSION	We use innovative science and engineering to anticipate and solve the most challenging national security problems
OBJECTIVE	In 10 years, we will have unleashed high-velocity engineering to counter global threats







The National Nuclear Security Administration (NNSA) Labs emphasize **trusted artificial intelligence (AI)** as a necessity to meet national security mission delivery.

The National Nuclear Security Administration (NNSA) Labs emphasize **trusted artificial intelligence (AI)** as a necessity to meet national security mission delivery.

- Recognize the value AI could provide to achieve mission goals
- Evaluating the credibility of current techniques poses challenges that may hinder widespread acceptance and use
- Sandia's mission needs set us apart from industry and academia due to reasons such as high-consequence applications



The National Nuclear Security Administration (NNSA) Labs emphasize **trusted artificial intelligence (AI)** as a necessity to meet national security mission delivery.

- Recognize the value AI could provide to achieve mission goals
- Evaluating the credibility of current techniques poses challenges that may hinder widespread acceptance and use
- Sandia's mission needs set us apart from industry and academia due to reasons such as high-consequence applications



The NNSA Labs must strike a balance between leveraging the advantages of ML while ensuring its responsible use for national security purposes.

G

MACHINE LEARNING AT SANDIA: EXAMPLE PROJECTS





Characterizing Variable Pathways Associated with a Volcanic Eruption



Image source: https://pubs.usgs.gov/fs/1997/fs113-97/







Trust defines the state of the decision maker.

Trustworthy defines the state of the model.

Credibility defines the technical basis of the model.

An example of a distinction in trust AI terminology at Sandia.





Trust defines the state of the decision maker.

 Human studies investigating how transparency and interactivity affect users' trust and performance with an algorithm

Trustworthy defines the state of the model.

Credibility defines the technical basis of the model.

EXPERIMENT

Participants interacted with an AI algorithm with differing levels of transparency and interactivity.

Four Possible GUIs

SLID is an algorithm which detects seismic arrival times and has two key parameters: smoothness and window size.





Trust defines the state of the decision maker.

 Human studies investigating how transparency and interactivity affect users' trust and performance with an algorithm

Trustworthy defines the state of the model.

• Projects considering counter-adversarial ML

Credibility defines the technical basis of the model.





Trust defines the state of the decision maker.

 Human studies investigating how transparency and interactivity affect users' trust and performance with an algorithm

Trustworthy defines the state of the model.

• Projects considering counter-adversarial ML

Credibility defines the technical basis of the model.

• Investigations into uncertainty quantification methods for ML and technical evaluations of explainability methods





MACHINE LEARNING AT SANDIA: ELEMENTS OF TRANSPARENCY

One area emphasized at Sandia is AI model transparency, which includes model interpretability and explainability

G

MACHINE LEARNING AT SANDIA: ELEMENTS OF TRANSPARENCY

One area emphasized at Sandia is AI model transparency, which includes model interpretability and explainability

Interpretability

Ability to directly use model to understand how algorithm makes decisions

$$\hat{y} = {\hat{eta}}_0 + {\hat{eta}}_1 x_1 + \dots + {\hat{eta}}_p x_p$$

Using interpretable model or adjusting black box models to contain interpretable parameters

MACHINE LEARNING AT SANDIA: ELEMENTS OF TRANSPARENCY



One area emphasized at Sandia is AI model transparency, which includes model interpretability and explainability

Interpretability

Ability to directly use model to understand how algorithm makes decisions

$$\hat{y} = {\hat{eta}}_0 + {\hat{eta}}_1 x_1 + \dots + {\hat{eta}}_p x_p$$

Explainability

Ability to indirectly use model to understand how algorithm makes decisions



Using interpretable model or adjusting black box models to contain interpretable parameters

Often post-hoc techniques Figure from LIME paper (<u>Ribeiro 2016</u>)

FORESTR

TOPOLOGY PATTERNS IN RANDOM FORESTS

MOTIVATION



MOTIVATION: RANDOM FORESTS

Random forests have become a common model used/considered in applications across Sandia

- Ensemble of decision/regression trees
- Introduces randomness in two places:
 - Each tree trained on a bootstrap sample of training data
 - Each split considers a random subset of features
- Advantages over a single tree:
 - Helps prevent overfitting
 - More robust to small variations in data





MOTIVATION: EXAMPLE DATA

Palmer Penguins

- Data: 342 penguins from Palmer Archipelago in Antarctica
- Goal: Predict species (Adelie, chinstrap, or gentoo)
 - Given four body measurements (bill length, bill depth, flipper length, body mass)
- Model: Random forests with 50 trees (max depth of 4 + full depth)





MOTIVATION: INTERPRETING TREE MODELS Single tree Typically considered interpretable • bill length mm \leq 41.6 samples = 153value = [89, 47, 97] class = Gentoobill depth mm \leq 14.9 flipper length mm ≤ 205.5 samples = 57samples = 96value = [82, 1, 1] value = [7, 46, 96] class = Gentooclass = Adelie



MOTIVATION: INTERPRETING TREE MODELS

- Single tree Typically considered interpretable
- Ensemble of trees Naturally becomes more difficult to interpret due to cognitive load



PREVIOUS WORK



Various approaches have been developed for gaining insight into how random forests use data for predictions

Breiman, L. Random Forests. Machine Learning 45, 5-32 (2001). https://doi.org/10.1023/A:1010933404324

PREVIOUS WORK: EXPLAINING RANDOM FORESTS

Various approaches have been developed for gaining insight into how random forests use data for predictions

Permutation Feature Importance (Breiman 2001)

Permute an input variable and quantify change in model performance





• Model: f

- Data: **X** with *n* obs and *p* variables
- Variables: X₁, X₂,..., X_p columns of
 X
- Metric: *m* computed with **X** and *f* (s.t. larger indicates better performance)

Procedure:

For variable $j \in \{1, \ldots, p\}$ and repetition $k \in \{1, \ldots, K\}$:

1. Create $ilde{X}_{j,k}$ by randomly permuting X_j

2. Create $ilde{\mathbf{X}}_{j,k}$ by replacing X_j with $ilde{X}_{j,k}$ in \mathbf{X}

```
3. Compute m_{j,k} with \tilde{\mathbf{X}}_{j,k} and f
```

PFI for j: Average change in model performance when j is randomly permuted

$$\mathscr{I}_j = m - rac{1}{K}\sum_{k=1}^K m_{j,k} = rac{1}{K}\sum_{k=1}^K \left(m - m_{j,k}
ight) = rac{1}{K}\sum_{k=1}^K \mathscr{I}_{j,k}$$

Various approaches have been developed for gaining insight into how random forests use data for predictions

Trace Plots (Urbanek 2010)

Designed to compare

- 1. variables used for splitting,
- 2. location of split points, and
- 3. hierarchical structure

Urbanek, S. (2008). Visualizing Trees and Forests. In: Handbook of Data Visualization. Springer Handbooks Comp.Statistics. Springer, Berlin, Heidelberg. <u>https://doi.org/</u> <u>10.1007/978-3-540-33037-0_11</u>







Various approaches have been developed for gaining insight into how random forests use data for predictions

Trace Plots (Urbanek 2010)

Designed to compare

- 1. variables used for splitting,
- 2. location of split points, and
- 3. hierarchical structure

Urbanek, S. (2008). Visualizing Trees and Forests. In: Handbook of Data Visualization. Springer Handbooks Comp.Statistics. Springer, Berlin, Heidelberg. <u>https://doi.org/</u> <u>10.1007/978-3-540-33037-0_11</u>





Split variable (ordered by random forest importance from left to right)

Various approaches have been developed for gaining insight into how random forests use data for predictions

Representative tree

- Identify a tree that is representative of the forest
- One approach: Find tree that has smallest average distance to all other trees

Clusters of trees

- Compute distances between trees
- Identify clusters via MDS, K-means, etc.

Banerjee, M., Y. Ding, and A. Noone (2012). "Identifying representative trees from ensembles". In: *Statistics in Medicine* 31.15, pp. 1601-1616. ISSN: 1097-0258. <u>10.1002/sim.4492</u>. Chipman, H. A., E. I. George, and R. E. McCulloch (1998). "Making sense of a forest of trees". In: *Proceedings of the 30th Symposium on the Interface*., pp. 84-92. <u>http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.42.2598</u>. Shannon, W.D. and Banks, D. (1999), Combining classification trees using MLE. Statist. Med., 18: 727-740. <u>https://doi.org/10.1002/(SICI)1097-0258(19990330)18:6<727::AID-SIM61>3.0.CO;2-2</u> Sies, A. and I. V. Mechelen (2020). "C443: a Methodology to See a Forest for the Trees". In: *Journal of Classification* 37.3, pp. 730-753. ISSN: 0176-4268. <u>https://link.springer.com/article/10.1007/s00357-019-09350-4</u>. Weinberg, A. I. and M. Last (2019). "Selecting a representative decision tree from an ensemble of decision-tree models for fast big data classification". In: *Journal of Big Data* 6.1, p. 23. <u>10.1186/s40537-019-0186-3</u>.

PREVIOUS WORK: REPRESENTATIVE AND CLUSTERING OF TREES



PREVIOUS WORK: REPRESENTATIVE AND CLUSTERING OF TREES

Identifying a representative tree or clusters of trees both require a distance metric

Banerjee, M., Y. Ding, and A. Noone (2012). "Identifying representative trees from ensembles". In: *Statistics in Medicine* 31.15, pp. 1601-1616. ISSN: 1097-0258. <u>10.1002/sim.4492</u>. Chipman, H. A., E. I. George, and R. E. McCulloch (1998). "Making sense of a forest of trees". In: *Proceedings of the 30th Symposium on the Interface*., pp. 84-92. <u>http://citeseerx.ist.psu.edu/viewdoc/summary?</u> doi=10.1.1.42.2598.

G

Banerjee, M., Y. Ding, and A. Noone (2012). "Identifying representative trees from ensembles". In: *Statistics in Medicine* 31.15, pp. 1601-1616. ISSN: 1097-0258. <u>10.1002/sim.4492</u>. Chipman, H. A., E. I. George, and R. E. McCulloch (1998). "Making sense of a forest of trees". In: *Proceedings of the 30th Symposium on the Interface*., pp. 84-92. <u>http://citeseerx.ist.psu.edu/viewdoc/summary?</u> <u>doi=10.1.1.42.2598</u>.

PREVIOUS WORK: REPRESENTATIVE AND CLUSTERING OF TREES

Identifying a representative tree or clusters of trees both require a distance metric Two approaches:

Comparing Predictions

Comparing Topology





և մ

PREVIOUS WORK: REPRESENTATIVE AND CLUSTERING OF TREES

Identifying a representative tree or clusters of trees both require a distance metric Two approaches:

Comparing Predictions

Comparing Topology

Fit metric Compares predictions from two trees (Chipman, George, and McCulloch (1998)

$${d_{FM}}\left({{T_1},{T_2}}
ight) = rac{1}{n}\sum\limits_{i = 1}^n {m\left({{\hat y_{i1}},{\hat y_{i2}}}
ight)}$$

where T_t is tree t with $t \in \{1, 2\}$, y_i is response, \hat{y}_{it} is fitted value, and m is a metric such as

- Regression metric: $m\left(\hat{y}_{i1},\hat{y}_{i2}
 ight)=\left(\hat{y}_{i1}-\hat{y}_{i2}
 ight)^2$
- Classification metric: $m\left(\hat{y}_{i1},\hat{y}_{i2}
 ight) = egin{cases} 1 & ext{if} \;\; \hat{y}_{i1}
 eq \hat{y}_{i2} \ 0 & ext{o.w.} \end{cases}$

Banerjee, M., Y. Ding, and A. Noone (2012). "Identifying representative trees from ensembles". In: Statistics in Medicine 31.15, pp. 1601-1616. ISSN: 1097-0258. 10.1002/sim.4492. Chipman, H. A., E. I. George, and R. E. McCulloch (1998). "Making sense of a forest of trees". In: Proceedings of the 30th Symposium on the Interface., pp. 84-92. http://citeseerx.ist.psu.edu/viewdoc/summary? doi=10.1.1.42.2598.

և մ

PREVIOUS WORK: REPRESENTATIVE AND CLUSTERING OF TREES

Identifying a representative tree or clusters of trees both require a distance metric Two approaches:

Comparing Predictions

Fit metric Compares predictions from two trees (Chipman, George, and McCulloch (1998)

 $d_{FM}\left(T_{1},T_{2}
ight)=rac{1}{n}\sum_{i=1}^{n}m\left(\hat{y}_{i1},\hat{y}_{i2}
ight)$

Comparing Topology

Covariate metric Compares split variables from two trees (Banerjee, Ding, and Noone (2012))

 $d_{CM}(T_1,T_2) = rac{ ext{Number of covariate mismatches for } T_1 ext{ and } T_2}{k}.$

where T_t is tree t with $t \in \{1,2\}$, y_i is response, \hat{y}_{it} is fitted value, and m is a metric such as

- Regression metric: $m\left(\hat{y}_{i1},\hat{y}_{i2}
ight)=\left(\hat{y}_{i1}-\hat{y}_{i2}
ight)^2$

• Classification metric:
$$m\left(\hat{y}_{i1},\hat{y}_{i2}
ight) = egin{cases} 1 & ext{if} \;\; \hat{y}_{i1}
eq \hat{y}_{i2} \ 0 & ext{o.w.} \end{cases}$$

Banerjee, M., Y. Ding, and A. Noone (2012). "Identifying representative trees from ensembles". In: *Statistics in Medicine* 31.15, pp. 1601-1616. ISSN: 1097-0258. <u>10.1002/sim.4492</u>. Chipman, H. A., E. I. George, and R. E. McCulloch (1998). "Making sense of a forest of trees". In: *Proceedings of the 30th Symposium on the Interface*., pp. 84-92. <u>http://citeseerx.ist.psu.edu/viewdoc/summary?</u> <u>doi=10.1.1.42.2598</u>.

PREVIOUS WORK: REPRESENTATIVE AND CLUSTERING OF TREES

Example of random forest tree clusters and corresponding central trees



Clusters identified using multidimensional scaling with fit metric and representative trees determined by tree with smallest average fit metric distance to all other trees in cluster.





FORESTR: PROJECT GOALS

FORESTR: Finding, Organizing, Representing, Explaining, Summarizing, and Thinning Random forests

Overview

- Use graph topology distance metric to identify tree topology patters patterns in forests
- Would like a distance metric for comparing tree topologies that is a proper mathematical distance to allow for the computation of summary statistics
- Future work could consider using metric to create ensemble with a reduced number of trees







Guo, X., Srivastava, A. & Sarkar, S. A Quotient Space Formulation for Generative Statistical Analysis of Graphical Data. *J Math Imaging Vis* 63, 735–752 (2021). <u>https://doi.org/10.1007/s10851-021-01027-1</u> Bal, Aditi Basu, et al. "Statistical Shape Analysis of Shape Graphs with Applications to Retinal Blood-Vessel Networks." *arXiv preprint arXiv:2211.15514* (2022). 67)

Adapt graph metric from Guo (2021) to trees, which has nice properties:

Proper mathematical distance, invariant to transformations, registration between graphs, and computes topological evolutions between graphs



Guo, X., Srivastava, A. & Sarkar, S. A Quotient Space Formulation for Generative Statistical Analysis of Graphical Data. *J Math Imaging Vis* 63, 735–752 (2021). <u>https://doi.org/10.1007/s10851-021-01027-1</u> Bal, Aditi Basu, et al. "Statistical Shape Analysis of Shape Graphs with Applications to Retinal Blood-Vessel Networks." *arXiv preprint arXiv:2211.15514* (2022).

FORESTR: TREE DISTANCE METRIC

For an ensemble of trees $t=1,\ldots,T$, let

- G_t = weighted graph associated with tree t
- (V_t, e_t) = nodes and edges associated with G_t



FORESTR: TREE DISTANCE METRIC

For an ensemble of trees $t=1,\ldots,T$, let

- G_t = weighted graph associated with tree t
- (V_t,e_t) = nodes and edges associated with G_t
- A_t = adjacency matrix associated with tree t



bill length mm \leq 41.6 samples = 153 value = [89, 47, 97] class = Gentoo flipper_length_mm <= 205.5 bill_depth_mm <= 14.9 samples = 57samples = 96 value = [82, 1, 1] value = [7, 46, 96] class = Gentoo class = Adelie bill_depth_mm <= 16.65 $body_mass_g \le 4100.0$ bill depth mm ≤ 17.65 samples = 1samples = 56samples = 34 samples = 62value = [0, 0, 1] value = [6.0, 41.0, 0.0] value = [1, 5, 96] value = [82, 1, 0] class = Gentoo class = Adelie class = Chinstrap class = Gentoo samples = 3samples = 53 samples = 29 samples = 5samples = 58samples = 4value = [2, 1, 0]value = [80, 0, 0]value = [2, 39, 0]value = [4, 2, 0]value = [0, 0, 96]value = [1, 5, 0]class = Adelie class = Adelie class = Chinstrap class = Adelie class = Gentoo class = Chinstrap

A tree in penguin random forest

Corresponding adjacency matrix

FORESTR: TREE DISTANCE METRIC

For an ensemble of trees $t=1,\ldots,T$, let

- G_t = weighted graph associated with tree t
- (V_t, e_t) = nodes and edges associated with G_t
- A_t = adjacency matrix associated with tree t





A tree in penguin random forest

GOAL For all pairs of trees t and u, want to compute distance between trees:

 $d_g([A_t], [A_u])$

Compute the distance between trees G_t and G_u as

$$d_g([A_t],[A_u]) = \min_{P\in\mathcal{P}} \|PA_tP^T - A_u\|^2 + \lambda \mathrm{Tr}(PD_{t,u})$$



Compute the distance between trees G_t and G_u as



where the optimization is implemented using Umeyama algorithm or fast approximate quadratic programming and

- *P* = permutation matrix
- \mathcal{P} = set of all permutation matrices of dimension n imes n
- λ = tuning parameter specifying weight to place on attributes
- $D_{t,u}$ = distance between node attributes of graphs t and u



Compute the distance between trees G_t and G_u as



where the optimization is implemented using Umeyama algorithm or fast approximate quadratic programming and

- *P* = permutation matrix
- \mathcal{P} = set of all permutation matrices of dimension n imes n
- λ = tuning parameter specifying weight to place on attributes
- $D_{t,u}$ = distance between node attributes of graphs t and u

For trees G_t and G_u ,

$$D_{t,u} = [d_{ij} = d(lpha_t(v_i^t), lpha_u(v_i^u))] \in \mathbb{R}^{n imes n},$$

where $\alpha_t(v_i^t)$ represents a node attribute and d is some distance metric.

Compute the distance between trees G_t and G_u as



where the optimization is implemented using Umeyama algorithm or fast approximate quadratic programming and

- *P* = permutation matrix
- \mathcal{P} = set of all permutation matrices of dimension n imes n
- λ = tuning parameter specifying weight to place on attributes
- $D_{t,u}$ = distance between node attributes of graphs t and u

For trees G_t and G_u ,

$$D_{t,u} = [d_{ij} = d(lpha_t(v_i^t), lpha_u(v_i^u))] \in \mathbb{R}^{n imes n},$$

where $\alpha_t(v_i^t)$ represents a node attribute and d is some distance metric.

To account for the hierarchical nature of trees, we consider node depth as an attribute.





Compute the distance between trees G_t and G_u as



where the optimization is implemented using Umeyama algorithm or fast approximate quadratic programming and

- *P* = permutation matrix
- \mathcal{P} = set of all permutation matrices of dimension n imes n
- λ = tuning parameter specifying weight to place on attributes
- $D_{t,u}$ = distance between node attributes of graphs t and u

For trees G_t and G_u ,

$$D_{t,u} = [d_{ij} = (lpha_t(v_i^t) - lpha_u(v_i^u))^2] \in \mathbb{R}^{n imes n},$$

where $\alpha_t(v_i^t)$ represents the node depth associated with node i in tree t.

To account for the hierarchical nature of trees, we consider node depth as an attribute.





Compute the distance between trees G_t and G_u as





Compute the distance between trees G_t and G_u as





Compute the distance between trees G_t and G_u as









Given a set of T graphs, denoted as $G_i \in \mathcal{G}, t = 1, \ldots, m$, and their respective adjacency matrices, $A_t \in \mathbb{R}^{n \times n}$, the mean graph is obtained by minimizing the sum of squared distances:

$$[A_\mu] = rgmin_{A\in \mathbb{R}^{n imes n}} \sum_{t=1}^T d_g([A],[A_i])^2$$

The mean graph is guaranteed to be unique, and the optimization is performed using a greedy optimization algorithm.



Given a set of T graphs, denoted as $G_i \in \mathcal{G}, t = 1, \ldots, m$, and their respective adjacency matrices, $A_t \in \mathbb{R}^{n \times n}$, the mean graph is obtained by minimizing the sum of squared distances:

$$[A_\mu] = rgmin_{A \in \mathbb{R}^{n imes n}} \sum_{t=1}^T d_g([A], [A_i])^2$$

The mean graph is guaranteed to be unique, and the optimization is performed using a greedy optimization algorithm.

This optimization results in violations of the properties of trees:



Given a set of T graphs, denoted as $G_i \in \mathcal{G}, t = 1, \ldots, m$, and their respective adjacency matrices, $A_t \in \mathbb{R}^{n \times n}$, the mean graph is obtained by minimizing the sum of squared distances:

$$[A_\mu] = rgmin_{A \in \mathbb{R}^{n imes n}} \sum_{t=1}^T d_g([A], [A_i])^2$$

The mean graph is guaranteed to be unique, and the optimization is performed using a greedy optimization algorithm.

This optimization results in violations of the properties of trees:

Currently, we compute a **central tree** as the three with the smallest distance to the mean graph. That is, let A_{μ} represent the mean graph, the central tree of a set of trees is computed as

$$A_{central} = rgmin_{A\in\{A_1,...,A_T\}} d_g(A,A_\mu).$$



Given a set of T graphs, denoted as $G_i \in \mathcal{G}, t = 1, \ldots, m$, and their respective adjacency matrices, $A_t \in \mathbb{R}^{n \times n}$, the mean graph is obtained by minimizing the sum of squared distances:

$$[A_\mu] = rgmin_{A \in \mathbb{R}^{n imes n}} \sum_{t=1}^T d_g([A], [A_i])^2$$

The mean graph is guaranteed to be unique, and the optimization is performed using a greedy optimization algorithm.

This optimization results in violations of the properties of trees:

Currently, we compute a **central tree** as the three with the smallest distance to the mean graph. That is, let A_{μ} represent the mean graph, the central tree of a set of trees is computed as

$$A_{central} = rgmin_{A\in\{A_1,...,A_T\}} d_g(A,A_\mu).$$





FORESTR: CENTRAL TREE WITHIN CLUSTERS



Central trees from each cluster

FORESTR: PRODUCT INSPECTION APPLICATION



FORESTR: PRODUCT INSPECTION APPLICATION

s500 data

- Information from a product inspection application (Kegelmeyer 2015)
- 20 numeric features associated with 1000 products identified as good or bad
- Separated into training and testing sets with 500 observations in each set

Predict good/bad

- 25 trees
- Maximum depth 5
- Fit using scikit-learn and all other default parameters



Parallel coordinate plots of s500 training (top) and testing (bottom) data. Lines are colored by the inspection labels, and the features are ordered by Gini importance.

Truth — -1 — 1

P. Kegelmeyer, T. M. Shead, J. Crussell, K. Rodhouse, D. Robinson, C. Johnson, D. Zage, W. Davis, J. Wendt, T. Cayton J. Doak, R. Colbaugh, K. Glass, B. Jones, and J. Shelburg. Counter adversarial data analytics Technical report, Sandia National Labs (SAND2015-3711), 2015.

6)

FORESTR: PRODUCT INSPECTION APPLIC









Dendrogram of complete linkage clusters based on pairwise distances

FORESTR: PRODUCT INSPECTION APPLICATION

Trees by Cluster

The number of trees in clusters 3 and 4 is much less than the other clusters (e.g., only 1 tree in cluster 3).

- Could these be outlier trees?
- How would removing them from the model affect model performance?



Trace plots of the tree clusters from the s500 random forest.

FORESTR: PRODUCT INSPECTION APPLICATION

The central trees help identify differences between the clusters:

- the depths at which leaf nodes begin to occur and
- the number of leaf nodes at a depth.



Central trees from each cluster

FORESTR: FUTURE RESEARCH DIRECTIONS



FORESTR: FUTURE RESEARCH DIRECTIONS

- Account for additional node attributes
 - (e.g., split feature)
- Consider approaches to identify repeating sub-graphs
 - (i.e., split patterns that are common among the forest)
- Use distances to reduce ensemble size and consider affect on predictive performance
 - Build up a model using a central tree or central trees from identified clusters
 - Trim down a model based on their distance to central tree or cluster central trees





67)

THANK YOU

KJGOODE@SANDIA.GOV GOODEKAT.GITHUB.IO

Technical Report Goode, Katherine Jean and Tucker, James Derek. "FORESTR: Finding, Organizing, Representing, Explaining, Summarizing, and Thinning Random forests." Sep. 2024. https://doi.org/10.2172/2472741.

