

# TESTS FOR LARGE TREE-STRUCTURED DATA

Karthik Bharath

School of Mathematical Sciences  
University of Nottingham

## IN THIS TALK..

Tree: acyclic graph with a 'root' which can be embedded on a plane.

## IN THIS TALK..

Tree: acyclic graph with a 'root' which can be embedded on a plane.

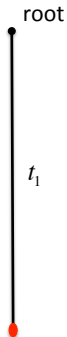
- A simple model for binary trees using Poisson process;
- Extend to general classes of **large** trees based on the **Continuum Random Tree**;
- Goodness-of-fit tests;
- Application to detecting brain tumor heterogeneity from images (Time permitting).

## MOTIVATION: MODEL FOR BINARY TREES

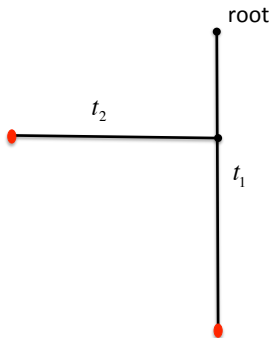
- Consider a non-homogeneous Poisson process with rate  $\lambda(t) = t$ .
- Let  $t_1, t_2, \dots$ , be inter-event times.

## MOTIVATION: MODEL FOR BINARY TREES

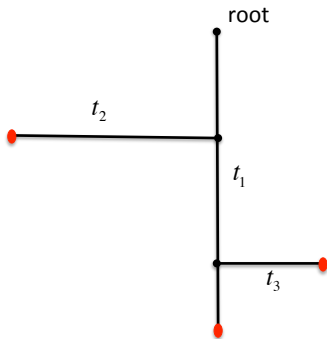
- Consider a non-homogeneous Poisson process with rate  $\lambda(t) = t$ .
- Let  $t_1, t_2, \dots$ , be inter-event times.



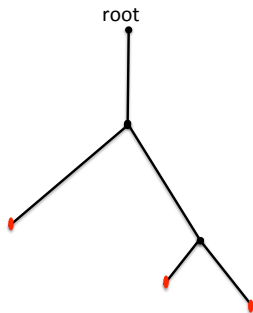
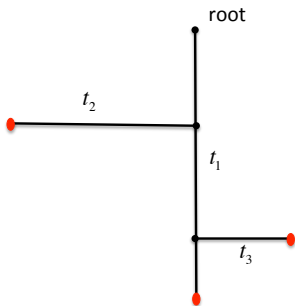
## MOTIVATION: MODEL FOR BINARY TREES



## MOTIVATION: MODEL FOR BINARY TREES



## MOTIVATION: MODEL FOR BINARY TREES



$$t_1 + t_2 + t_3 = \text{total path length of binary tree}$$



## MOTIVATION: MODEL FOR BINARY TREES

With  $n$  inter-event times, a binary tree  $\tau(n)$  with  $n$  leaves or terminal nodes,  $2n$  vertices and  $2n - 1$  edges is constructed.

## MOTIVATION: MODEL FOR BINARY TREES

With  $n$  inter-event times, a binary tree  $\tau(n)$  with  $n$  leaves or terminal nodes,  $2n$  vertices and  $2n - 1$  edges is constructed.

### PROPOSITION

*From the properties of the Poisson process with rate  $t$ ,  $\tau(n)$  can be given the density*

$$f(\tau(n)) = \left[ \prod_{i=1}^{n-1} \frac{1}{2i-1} \right]^{-1} \frac{1}{2^{n-1}} s e^{-s^2/2}, \quad s = \sum_{i=1}^{2n-1} t_i.$$

- $f(\cdot)$  is exchangeable with respect to the branch lengths.
- $f(\cdot)$  is independent of the “shape” of the tree.
- $f(\cdot)$  is ‘consistent’: removal of a leaf from  $\tau(n)$  results in a tree with density  $f(\tau(n-1))$ .

## TOTAL PATH LENGTH

### PROPOSITION

*Suppose  $\tau(k)$  is a binary tree with  $k$  leaves and branch lengths  $x_1, \dots, x_{2k-1}$  generated from the non-homogeneous Poisson model. Then the total path length  $\sum_{i=1}^{2k-1} x_i$  follows a Gamma distribution with shape  $k$  and scale 2.*

## GoF TEST FOR BINARY TREES: ONE-SAMPLE

Suppose  $\tau(\mathbf{n}) = (\tau(n_1), \dots, \tau(n_p))$  is an independent sample of binary trees from  $\pi_\tau$ .

### THEOREM

Consider the critical function

$$\phi(\mathbf{n}, \alpha) = \begin{cases} 1 & \text{if } \sum_{i=1}^p s_i > \chi_{1-\alpha, 2\sum_{i=1}^p n_i} \\ 0 & \text{otherwise,} \end{cases}$$

where  $s_i$  is the sum of the branch lengths of  $\tau(n_i)$  and  $\chi_{\alpha, b}$  denotes the  $\alpha$ th percentile of a Chi-square distribution with  $b$  degrees of freedom. For the hypotheses  $H_0 : \pi_\tau = f$  against  $H_1 : \pi_\tau \neq f$ , where  $f$  is the density from the non-homogeneous Poisson model,  $E_{H_0} \phi(\mathbf{n}, \alpha) = \alpha$ .

## GOF TEST FOR BINARY TREES: TWO-SAMPLE

Suppose  $\boldsymbol{\tau}(\mathbf{n}) = (\tau(n_1), \dots, \tau(n_p))$  and  $\boldsymbol{\eta}(\mathbf{m}) = (\eta(m_1), \dots, \eta(m_q))$  are independent samples of binary trees from  $\pi_\tau$  and  $\pi_\eta$  respectively.

### THEOREM

Let  $r_j$  denote the sum of the branch lengths of  $\eta(m_j)$ , and without loss of generality assume that  $\sum_{i=1}^p s_i > \sum_{j=1}^q r_j$ . Then, the critical function

$$\psi(\mathbf{n}, \mathbf{m}, \alpha) = \begin{cases} 1 & \text{if } \frac{\sum_{i=1}^p s_i}{\sum_{j=1}^q r_j} > \left( \frac{\sum_{i=1}^p n_i}{\sum_{j=1}^q m_j} \right) F_{1-\alpha, 2 \sum_{i=1}^p n_i, 2 \sum_{j=1}^q m_j}; \\ 0 & \text{otherwise,} \end{cases}$$

where  $F_{\alpha, a, b}$  is the  $\alpha$ th percentile of an  $F$  distribution with  $a$  and  $b$  degrees of freedom, for testing  $H_0 : \pi_\tau = \pi_\eta = f$ , is such that  $E_{H_0} \psi(\mathbf{n}, \mathbf{m}, \alpha) = \alpha$ .

## NON-BINARY TREES

- Can the model and tests for binary trees be extended to other types of trees?

## NON-BINARY TREES

- Can the model and tests for binary trees be extended to other types of trees?

Yes, in an asymptotic sense for 'large' trees with some modifications.

## MOTIVATION

“ Given a sequence of discrete combinatorial random structures of size  $n$ , does there exist a continuous structure representing their  $n \rightarrow \infty$  limit? If so, then for many questions about the size- $n$  object one can obtain the  $n \rightarrow \infty$  limit by simply asking the same question of the limit structure. This is the *weak convergence paradigm*. ”

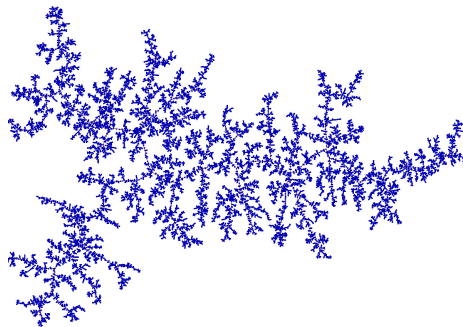
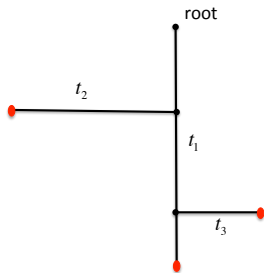
– *David Aldous*.

For random trees, the continuous structure is termed the **Continuum Random Tree** (CRT).



# CONTINUUM RANDOM TREE (CRT)

As  $n \rightarrow \infty$ ,



(<http://www.normalesup.org/kortchem/english.html>)

## WHAT IS THE CRT?

- It can thought of as the closure of the union of binary trees  $\cup_n \tau(n)$ .
- From a stochastic process perspective, its 'finite-dimensional distributions' are the the binary trees studied earlier!

## WHAT IS THE CRT?

- It can thought of as the closure of the union of binary trees  $\cup_n \tau(n)$ .
- From a stochastic process perspective, its 'finite-dimensional distributions' are the the binary trees studied earlier!

It is realized as:

- the weak limit, as vertices grow without bound, of Galton–Watson process conditioned on total progeny;
- an object constructed from the dense set of local minima of a standard Brownian excursion process.

## General tree models with the CRT

## REPRESENTATION OF TREE

- We will use a more intuitive representation of a tree  $\tau_n$  with  $n$  vertices and  $n - 1$  edges:

$$\tau_n = (\mathcal{V}(\tau_n), \mathcal{E}(\tau_n)),$$

where  $\mathcal{V}(\tau_n) = (\text{root}, v_1, \dots, v_{n-1})$  is the vertex-set and  $\mathcal{E}(\tau_n) = (e_1, \dots, e_{n-1})$  is the edge-set.

- In other words,  $\tau_n$  is a point in  $\mathcal{T}_n \times \mathbb{R}_+^{n-1}$  where  $\mathcal{T}_n$  is the set of all combinatorial trees with  $n$  vertices.

## CONDITIONED GALTON-WATSON TREE MODELS

- Suppose  $\xi$  is non-negative integer-valued r.v. with distribution  $(\pi_k : k \geq 0)$ .
- Construct a tree  $\tau$  recursively starting with root and giving each node a number of children that is an independent copy of  $\xi$ . This induces a distribution on  $\tau$ :

$$P(\tau = t) = \prod_{v \in \mathcal{V}(t)} \pi_{o(v,t)},$$

where  $o(v, t)$  is the out-degree or the number of children of vertex  $v$  in tree  $t$ .

## CONDITIONED GALTON-WATSON TREE MODELS

- Suppose  $\xi$  is non-negative integer-valued r.v. with distribution  $(\pi_k : k \geq 0)$ .
- Construct a tree  $\tau$  recursively starting with root and giving each node a number of children that is an independent copy of  $\xi$ . This induces a distribution on  $\tau$ :

$$P(\tau = t) = \prod_{v \in \mathcal{V}(t)} \pi_{o(v,t)},$$

where  $o(v, t)$  is the out-degree or the number of children of vertex  $v$  in tree  $t$ .

- **Conditioned Galton-Watson** (CGW) trees are family trees of Galton-processes conditioned on total progeny. The distribution of a CGW tree  $\tau_n$  conditioned on  $n$  vertices is then

$$P(\tau_n = t) \propto \prod_{v \in \mathcal{V}(t)} \pi_{o(v,t)} \quad \text{on } \{t : \text{cardinality of } \mathcal{V}(t) = n\}.$$

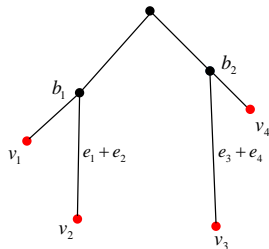
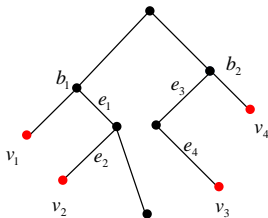
## CONDITIONED GALTON-WATSON TREE MODELS

- $\pi$  is referred to as the offspring distribution.
- CGW represent a broad class of trees which can serve as probability models.
  - ① *Plane trees* (ordered trees): CGW trees with offspring distribution given by a Geometric distribution with success probability  $1/2$ , and  $\sigma^2 = 2$ ;
  - ② *Binary trees*: CGW trees with vertices containing 0,1 or 2 children with a Binomial distribution with 2 trials and success probability  $1/2$ , and  $\sigma^2 = 1/2$ ;
  - ③ *Strict binary trees*: CGW trees with vertices containing either 0 or 2 children with equal probability  $1/2$ , and  $\sigma^2 = 1$ ;
  - ④ *Unary-binary trees*: CGW trees with vertices containing 0, 1 or 2 children each with probability  $1/3$ , and  $\sigma^2 = 2/3$ ;
  - ⑤ *m-ary trees*: CGW trees with vertices containing 0, 1,  $\dots$ ,  $m$  for  $m > 3$  children with distribution given by a Binomial with  $m$  trials and success probability  $1/m$ , and  $\sigma^2 = \frac{m-1}{m}$ .



## DISTRIBUTION OF CRT: LCA TREES

For a CGW tree  $\tau_n = (\mathcal{V}(\tau_n), \mathcal{E}(\tau_n))$ , consider the Least Common Ancestor (LCA) subtree spanned by a randomly chosen subset  $B$  of the **leaves**, including the root.



## DISTRIBUTION OF CRT: LCA SUBTREES (ALDOUS 1993)

For a CGW tree  $\tau_n = (\mathcal{V}(\tau_n), \mathcal{E}(\tau_n))$ , suppose  $|B|$  is  $k < n$  and denote the LCA tree spanned by  $B$  as  $LCA(\tau_n, B)$ . Then

*As  $n \rightarrow \infty$ , for a fixed  $k$ ,  $LCA(\tau_n, B)$  “converges weakly” to a binary tree,  $\mathcal{L}(k)$ , obtained via the non-homogeneous Poisson process model.*

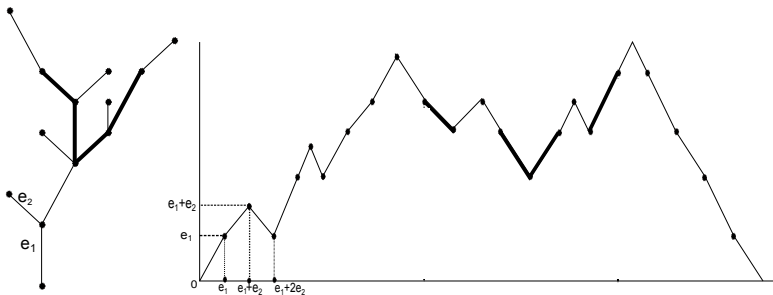
- $\mathcal{L}(k)$  can be viewed as “finite-dimensional projection” of CRT.
- Key point: Distribution of CRT is completely specified by distribution of  $\{\mathcal{L}(k), k \geq 1\}$ .

# DISTRIBUTION OF CRT: DYCK PATH REPRESENTATION

For a CGW tree with  $n$  vertices, define a continuous function  $H_n : [0, 2n] \rightarrow \mathbb{R}_{\geq 0}$  such that

$$H_n(s) = d(\text{root}, v),$$

where  $v$  is the vertex obtained during the depth-first walk such that the sum of the edges traversed till  $v$  is  $s$ .



## ALDOUS' REMARKABLE RESULT

### THEOREM

Let  $\tau_n$  be a CGW tree with offspring distribution with mean 1 and variance  $\sigma^2 \in (0, \infty)$ .

Let  $H_n(k)$ ,  $0 \leq k \leq 2n$  be the Dyck path associated with  $\tau_n$ . Then, as  $n \rightarrow \infty$ ,

$$\left\{ \frac{1}{\sqrt{n}} H_n([2nt]), 0 \leq t \leq 1 \right\} \Rightarrow \left\{ \frac{2}{\sigma} B_t^{\text{ex}} : 0 \leq t \leq 1 \right\}$$

where  $B^{\text{ex}}$  is the standard Brownian excursion.

## ALDOUS' EVEN MORE REMARKABLE RESULT

### THEOREM

*The Brownian excursion which arises as the limit of the normalized Dyck path of a conditioned Galton-Watson tree is the "Dyck path" of the CRT.*

# DISTANCE FROM ROOT OF RANDOMLY CHOSEN VERTEX

## PROPOSITION

Let  $U$  be uniform on  $[0, 1]$ . Consider the functional  $B^{\text{ex}}(U)$ . For a fixed  $\sigma^2$ , the class of distributions induced through  $U \mapsto \frac{2}{\sigma} B^{\text{ex}}(U)$  characterises the law of  $\frac{2}{\sigma} B^{\text{ex}}$ , with  $\frac{2}{\sigma} B^{\text{ex}}(U)$  following a Rayleigh distribution with scale  $1/\sigma$ .

# DISTANCE FROM ROOT OF RANDOMLY CHOSEN VERTEX

## PROPOSITION

Let  $U$  be uniform on  $[0, 1]$ . Consider the functional  $B^{\text{ex}}(U)$ . For a fixed  $\sigma^2$ , the class of distributions induced through  $U \mapsto \frac{2}{\sigma} B^{\text{ex}}(U)$  characterises the law of  $\frac{2}{\sigma} B^{\text{ex}}$ , with  $\frac{2}{\sigma} B^{\text{ex}}(U)$  following a Rayleigh distribution with scale  $1/\sigma$ .

## PROPOSITION

On an ordered conditioned Galton–Watson tree  $\tau_n$  with offspring variance  $\sigma^2$ , suppose  $v$  is a vertex chosen according to a uniform distribution on  $\mathcal{V}(\tau_n)$ . Then, the random variable

$$n^{-1/2} d(\text{root}, v) \xrightarrow{d} W,$$

as  $n \rightarrow \infty$ , where  $W$  is a Rayleigh distributed random variable with scale  $1/\sigma$ .

## GoF TESTS FOR CGW TREES

### LCA-based tests:

- For each tree in a sample, choose a subset of leaves at random, and construct LCA trees—each LCA tree will be a strict binary tree;
- From a non-homogeneous Poisson process with rate  $\lambda(t) = \sigma^2 t$ , obtain a parametric class of densities seen earlier for binary trees;
- With a consistent estimator for  $\sigma^2$ , the  $\chi^2$  and  $F$  GoF tests for binary trees are valid, as the number of vertices in each tree grow without bound.



## GoF TESTS FOR CGW TREES

### LCA-based tests:

- For each tree in a sample, choose a subset of leaves at random, and construct LCA trees—each LCA tree will be a strict binary tree;
- From a non-homogeneous Poisson process with rate  $\lambda(t) = \sigma^2 t$ , obtain a parametric class of densities seen earlier for binary trees;
- With a consistent estimator for  $\sigma^2$ , the  $\chi^2$  and  $F$  GoF tests for binary trees are valid, as the number of vertices in each tree grow without bound.

### Dyck path-based tests:

- For each tree in a sample, choose a vertex at random, and record its normalized distance from the root;
- Consider a  $\sigma^2$ -parameterized class of Rayleigh densities, and construct a consistent estimator of  $\sigma^2$ ;
- Noting that if  $X$  is Rayleigh distributed with scale  $1/b$ , then  $X^2$  is Chi-square distributed with 2 degrees of freedom scaled by  $1/b^2$ , we get similar GoF tests.

## PERFORMANCE OF LCA-BASED TEST

One-sample

Distribution	$n = 10$		$n = 100$		$n = 1000$	
	$\chi^2$	perm	$\chi^2$	perm	$\chi^2$	perm
Geo(0.5)	0.09	0.15	0.05	0.08	0.03	0.09
Bin(0.5)	0.13	0.08	0.04	0.03	0.01	0.01
Bin(0.35)	0.10	0.16	0.12	0.07	0.06	0.08
GW-Bin(2,0.5)	0.78	0.91	0.91	0.97	0.99	1.00
Phylo.bd	0.81	0.83	0.89	0.91	0.98	0.94
<b>Phylo.coal</b>	<b>0.26</b>	<b>0.37</b>	<b>0.14</b>	<b>0.21</b>	<b>0.11</b>	<b>0.08</b>

## Application: Detecting tumor heterogeneity with MRI

Joint work with collaborators at M.D. Anderson Cancer Center, Houston.

## TUMOUR HETEROGENEITY

- Variety of genetic, cellular and molecular mutations occur during the course of tumour development or during the course of a treatment.
- Classification of tumours using clustering of pixel intensities is popular.
- Typically histograms are compared with few parameters.

## TUMOUR HETEROGENEITY

- Variety of genetic, cellular and molecular mutations occur during the course of tumour development or during the course of a treatment.
- Classification of tumours using clustering of pixel intensities is popular.
- Typically histograms are compared with few parameters.

**Our approach:** Capture heterogeneity through variations in pixel intensities via hierarchical relationships between groups of pixels.

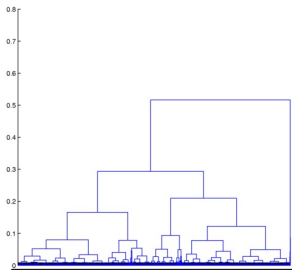
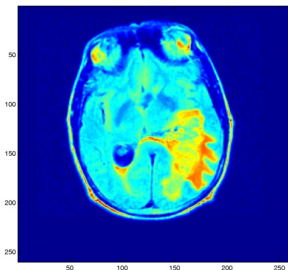
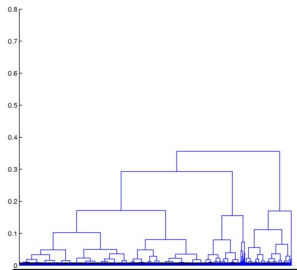
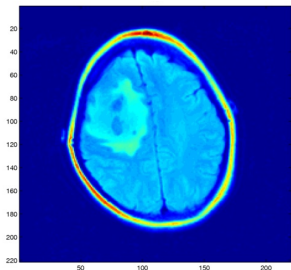
## DATASET

- Images of 82 patients with histologically confirmed GBM and molecular data from The Cancer Genome Atlas (TCGA) database (<https://www.cancerimagingarchive.net/>).
- T1-post and T2-FLAIR images were registered spatially followed by intensity bias correction using Medical Image Processing Analysis and Visualization software (v 6.0).13.
- The tumor region was segmented semi-automatically in 3D using the Medical Image Interaction Toolkit ([MITK.org](http://MITK.org)) with in-plane resolution of  $1\text{mm} \times 1\text{mm}$ .

## DATASET

- Images of 82 patients with histologically confirmed GBM and molecular data from The Cancer Genome Atlas (TCGA) database  
(<https://www.cancerimagingarchive.net/>).
- T1-post and T2-FLAIR images were registered spatially followed by intensity bias correction using Medical Image Processing Analysis and Visualization software (v 6.0).13.
- The tumor region was segmented semi-automatically in 3D using the Medical Image Interaction Toolkit ([MITK.org](http://MITK.org)) with in-plane resolution of  $1\text{mm} \times 1\text{mm}$ .
- T1 and T2 intensities from the segmented regions were grouped with **agglomerative hierarchical clustering to obtain dendrograms of image intensities.**

# DATASET





## DENDROGRAMS FROM HIERARCHICAL CLUSTERING

These dendrograms are **ultrametric** trees, and do not correspond to CGW trees. In fact, they can be generated by a coalescent process used in phylogenetics (recall poor power against phylo.coal).

## DENDROGRAMS FROM HIERARCHICAL CLUSTERING

These dendrograms are **ultrametric** trees, and do not correspond to CGW trees. In fact, they can be generated by a coalescent process used in phylogenetics (recall poor power against phylo.coal).

However, there is a surprising connection: if the leaves of an ultrametric binary tree,  $\tau(n)$  with  $n$  leaves, arising from a hierarchical clustering method are exchangeable in the sense that the distribution of the the leaves is invariant to permutation, then the *LCA*-trees converge in distribution, as  $n \rightarrow \infty$ , to the family  $\mathcal{L}(k)$  of subtrees which characterise the CRT. (Ph.D. dissertation of Chris Haulk [2012])

## TWO-SAMPLE TEST TO DETECT HETEROGENEITY

- Using the survival times, we created two groups of patients: those with survival times of utmost 12 months and those exceeding 12 months.
- The 12-month cut-off corresponded to a certain genetic classification— this was based on recommendations by neuroscientists.
- Differences in groups was detected by LCA-based test at 1% significance level.
- Naive Bayes classifier with the likelihood from LCA trees, provided 69% classification accuracy.

Details available in a paper on Arxiv: *Statistical Tests for Large Tree-structured Data*.

Details available in a paper on Arxiv: *Statistical Tests for Large Tree-structured Data*.

“If you can’t convince them, confuse them.”

—*Harry S. Truman*.