

APPENDIX

MATHEMATICAL DEFINITIONS

The trees we use are members of a normed vector space: we can add them, find out how far apart they are and interpolate between them. In principle, we can run clustering algorithms to find configurations that are similar, and identify when a model has left one cluster and entered another.

DEFINITION 1. Let \mathcal{S} be a set of labels, and let \mathbb{F} be a field like the real or complex numbers. Then we define a node \mathbf{n} as a triplet of the form (s, v, \mathbf{E}) , with $v \in \mathbb{F}$, $s \in \mathcal{S}$, and the set \mathbf{E} is a (possibly empty) set of nodes of the same form with the restriction that no two nodes in \mathbf{E} may have the same label in their first ordinate. We also define the triple $\mathbf{O} = (\emptyset, 0, \emptyset)$, which we will call the null tree, and define \mathbf{T} to be the union of the set of all trees composed of a finite number of these nodes.

The ordinates of $\mathbf{u} = (\mathbf{u}_s, \mathbf{u}_v, \mathbf{u}_E)$ in \mathbf{T} correspond to its label, value, and extension set. An element of \mathbf{u}_E will be called an extension.

In our discussion, it will help to have a few more descriptive terms.

A node with an empty extension set is called a *simple* node, if this node happens to be a root node, then it is a simple tree.

Two trees, $\mathbf{u}, \mathbf{v} \in \mathbf{T}$ are called *compatible* if either $\mathbf{u}_s = \mathbf{v}_s$ or at least one of \mathbf{u} and \mathbf{v} is \mathbf{O} .

We define the *depth* of a tree with:

$$\text{depth}(\mathbf{u}) = \begin{cases} 0 & \text{if } \mathbf{u} = (\emptyset, 0, \emptyset) = \mathbf{O} \\ 1 & \text{if } \mathbf{u} \text{ is a simple node} \\ 1 + \max(\{\text{depth}(\mathbf{v}) : \forall \mathbf{v} \in \mathbf{u}_E\}) & \text{otherwise.} \end{cases}$$

We will also define for $\mathbf{u} \in \mathbf{T}$,

$$\text{trim}(\mathbf{u}) = \begin{cases} \mathbf{O} & \text{if } \mathbf{u} = \mathbf{O} \\ \mathbf{O} & \text{if } \mathbf{u} \text{ is simple} \\ (\mathbf{u}_s, \mathbf{u}_v, \{\text{trim}(\mathbf{e}) : \forall \mathbf{e} \in \mathbf{u}_E\} \setminus \{\mathbf{O}\}) & \text{otherwise.} \end{cases}$$

The cardinality of a tree is the number of nodes it contains. Specifically,

$$\|\mathbf{u}\|_{\mathbf{T}} = \begin{cases} 0 & \text{if } \mathbf{u} = \mathbf{O} \\ 1 + \sum_{\mathbf{e} \in \mathbf{u}_E} \|\mathbf{e}\|_{\mathbf{T}} & \text{otherwise.} \end{cases}$$

Simple nodes are the only nodes that have a cardinality of one, and \mathbf{O} is the only node or tree with a cardinality of zero.

The support of a tree is the number of nodes which have a non-zero value.

$$\text{supp}(\mathbf{u}) = \begin{cases} 0 + \sum_{e \in \mathbf{u}_E} \text{supp}(e) & \text{if } \mathbf{u}_v = 0 \\ 1 + \sum_{e \in \mathbf{u}_E} \text{supp}(e) & \text{otherwise} \end{cases}.$$

Related is the fundamental support of a tree, which only counts nodes with no zero valued nodes in their connection to the root node

$$\text{fund}(\mathbf{u}) = \begin{cases} 0 & \text{if } \mathbf{u}_v = 0 \\ 1 + \sum_{e \in \mathbf{u}_E} \text{fund}(e) & \text{otherwise} \end{cases}.$$

Clearly the support of a tree, $\text{supp } \mathbf{u}$, must lie in the domain $[0, \|\mathbf{u}\|_{\mathcal{T}}]$ and $\text{fund } \mathbf{u} \leq \text{supp}(\mathbf{u})$.

The *overlap* between two trees is defined

$$\text{overlap}(\mathbf{u}, \mathbf{v}) = \begin{cases} 0 & \text{if } \mathbf{u} = \mathbf{O} \text{ or } \mathbf{v} = \mathbf{O} \text{ or } \mathbf{u}_s \neq \mathbf{v}_s \\ 1 + \sum_{\substack{e \in \mathbf{u}_E \\ f \in \mathbf{v}_E}} \text{overlap}(e, f) & \text{otherwise} \end{cases}$$

Clearly two trees, \mathbf{u} and \mathbf{v} , are compatible if and only if $\text{overlap}(\mathbf{u}, \mathbf{v}) \neq 0$; they will be said to *completely overlap* if $\|\mathbf{u}\|_{\mathcal{T}} = \|\mathbf{v}\|_{\mathcal{T}} = \text{overlap}(\mathbf{u}, \mathbf{v})$.

We can now define scalar multiplication, and tree addition.

DEFINITION 2. Take $a \in \mathbb{F}$ and $\mathbf{u} \in \mathcal{T}$, then

$$a\mathbf{u} = \begin{cases} \mathbf{O} & \text{if } \mathbf{u} = \mathbf{O} \\ (\mathbf{u}_s, a\mathbf{u}_v, \{af : f \in \mathbf{u}_E\} \setminus \{\mathbf{O}\}) & \text{otherwise.} \end{cases} \quad (1)$$

DEFINITION 3. Let \mathbf{u} and \mathbf{v} be compatible elements of \mathcal{T} . Then taking the symbol $+$ to be addition in the field \mathbb{F} , we extend it to addition in \mathcal{T} so that for nodes \mathbf{u} and \mathbf{v} ,

$$\mathbf{u} + \mathbf{v} = \begin{cases} \mathbf{O} & \text{if } \mathbf{u} = \mathbf{v} = \mathbf{O} \\ \mathbf{u} & \text{if } \mathbf{u} \neq \mathbf{O} \text{ and } \mathbf{v} = \mathbf{O} \\ \mathbf{v} & \text{if } \mathbf{u} = \mathbf{O} \text{ and } \mathbf{v} \neq \mathbf{O} \\ (\mathbf{u}_s, \mathbf{u}_v + \mathbf{v}_v, \emptyset) & \text{if } \mathbf{u}_E, \mathbf{v}_E = \emptyset \\ \left(\mathbf{u}_s, \mathbf{u}_v + \mathbf{v}_v, \left(\{f + g : f \in \mathbf{u}_E \text{ and } g \in \mathbf{v}_E \text{ and } f_s = g_s\} \right. \right. \\ \quad \cup \{f : f \in \mathbf{u}_E \text{ and } f_s \neq g_s \forall g \in \mathbf{v}\} \\ \quad \left. \left. \cup \{g : g \in \mathbf{v}_E \text{ and } g_s \neq f_s \forall f \in \mathbf{u}\} \right) \setminus \{\mathbf{O}\} \right) & \text{otherwise.} \end{cases} \quad (2)$$

DEFINITION 4. Let \mathbf{u} and \mathbf{v} be compatible elements of \mathcal{T} . Then we define inner-multiplication between the two nodes

$$\mathbf{u} \cdot \mathbf{v} = (\mathbf{u}_s, \mathbf{u}_v \mathbf{v}_v, \{f \cdot g : f \in \mathbf{u}_E, g \in \mathbf{v}_E \text{ and } f_s = g_s\} \setminus \{\mathbf{O}\}) \quad (3)$$

It can be shown that \mathcal{T} with scalar multiplication and tree addition forms a vector space. This isn't quite enough to give us distances between trees, however, so we define a semi-norm

DEFINITION 5. Let \mathbf{u} be an element of \mathcal{T} . Then we can define a semi-norm over \mathcal{T}

$$\langle \mathbf{u} \rangle = \begin{cases} 0 & \text{if } \mathbf{u} = \mathbf{O} \\ |\mathbf{u}_v| & \text{if } \mathbf{u}_E = \emptyset \\ |\mathbf{u}_v| + \sum_{e \in \mathbf{u}_E} \langle e \rangle & \text{otherwise.} \end{cases}$$

It is clear that $\langle \mathbf{u} \rangle$ will always be non-negative, and the only shortcoming is that we can have a node \mathbf{u} with $\langle \mathbf{u} \rangle = 0$, but $\mathbf{u} \neq \mathbf{O}$. In order to turn this into a normed vector space we take the set $\mathfrak{D} = \{\mathbf{u} : \mathbf{u} \in \mathcal{T} \text{ and } \langle \mathbf{u} \rangle = 0\}$ and we construct an equivalence relation on \mathcal{T} by the rule $[\mathbf{u}] \equiv [\mathbf{v}]$ if and only if there exist $\mathbf{z}_u, \mathbf{z}_v \in \mathfrak{D}$ such that $\mathbf{u} + \mathbf{z}_u = \mathbf{v} + \mathbf{z}_v$. It can be shown that scalar multiplication, tree addition, and the semi-norm behave appropriately with respect to the equivalence classes. This means that if we identify elements of \mathcal{T} with their equivalence class, then we can take $\langle \mathbf{u} \rangle$ to be a norm and that it induces a distance function

$$d(\mathbf{u}, \mathbf{v}) = \langle \mathbf{u} - \mathbf{v} \rangle$$

DEFINITION 6. We define the functions mask and $\overline{\text{mask}}$ that set the values associated with particular nodes in a tree to $v \in \mathbb{F}$. Specifically, if $\mathcal{L} \subset \mathcal{T}$, then

$$\text{mask}(\mathbf{u}, \mathcal{L}, v) = \begin{cases} (\mathbf{u}_s, v, \{\text{mask}(\mathbf{f}, \mathcal{L}, v) : \mathbf{f} \in \mathbf{u}_E\}) & \text{if } \mathbf{u}_s \in \mathcal{L} \\ (\mathbf{u}_s, \mathbf{u}_v, \{\text{mask}(\mathbf{f}, \mathcal{L}, v) : \mathbf{f} \in \mathbf{u}_E\}) & \text{otherwise} \end{cases} \quad (4)$$

and

$$\overline{\text{mask}}(\mathbf{u}, \mathcal{L}, v) = \begin{cases} (\mathbf{u}_s, v, \{\overline{\text{mask}}(\mathbf{f}, \mathcal{L}, v) : \mathbf{f} \in \mathbf{u}_E\}) & \text{if } \mathbf{u}_s \notin \mathcal{L} \\ (\mathbf{u}_s, \mathbf{u}_v, \{\overline{\text{mask}}(\mathbf{f}, \mathcal{L}, v) : \mathbf{f} \in \mathbf{u}_E\}) & \text{otherwise.} \end{cases} \quad (5)$$

$\text{mask}(\mathbf{u}, \mathcal{L}, 0)$ would return a tree similar to \mathbf{u} , but all its nodes that have labels in \mathcal{L} would have values of zero.

We also define several functions which prune or select a child from a tree's extension set. This function returns only the part of \mathbf{u} which overlaps \mathbf{p} ,

$$\text{excise}(\mathbf{u}, \mathbf{p}) = \begin{cases} (\mathbf{u}_s, \mathbf{u}_v, \{\text{excise}(\mathbf{f}, \mathbf{g}) : \mathbf{f} \in \mathbf{u}_E \text{ and } \mathbf{g} \in \mathbf{p} \text{ and } \mathbf{f}_s = \mathbf{g}_s\} \setminus \{\mathbf{O}\}) & \text{if } \mathbf{u}_s = \mathcal{L} \\ \mathbf{O} & \text{if } \mathbf{u}_s \neq \mathbf{p}_s \\ (\mathbf{u}_s, \mathbf{u}_v, \emptyset) & \text{otherwise,} \end{cases} \quad (6)$$

and this one either returns an appropriately labelled child from the extension set (a branch), or \mathbf{O} .

$$\text{child}(\mathbf{u}, \ell) = \begin{cases} \mathbf{f} & \text{if } \mathbf{f}_s = \ell \text{ and } \mathbf{f} \in \mathbf{u}_E \\ \mathbf{O} & \text{otherwise.} \end{cases} \quad (7)$$

We now finish with a definition of a function, $\Delta(\cdot)$, that gives us a measure of the degree of divergence between two trees.

DEFINITION 7. The degree of deviation between two trees, \mathbf{u} and \mathbf{v} is given by the expression

$$\Delta(\mathbf{u}, \mathbf{v}) = \begin{cases} \|\mathbf{u}\|_{\mathcal{T}} + \|\mathbf{v}\|_{\mathcal{T}} - 2\text{overlap}(\mathbf{u}, \mathbf{v}) + \langle \mathbf{u} - \mathbf{v} \rangle & \text{if } \mathbf{u} \text{ and } \mathbf{v} \text{ are compatible} \\ \|\mathbf{u}\|_{\mathcal{T}} + \|\mathbf{v}\|_{\mathcal{T}} & \text{otherwise.} \end{cases} \quad (8)$$

The rationale behind this definition is that if trees \mathbf{u} and \mathbf{v} are identical, then $\Delta(\mathbf{u}, \mathbf{v})$ will be zero. We also want nodes that aren't common to both trees to count as differences.