

---

# Hyperbolic Embeddings of Supervised Models

---

Anonymous Author(s)

Affiliation  
Address  
email

## Abstract

1 Models of hyperbolic geometry have been successfully used in ML for two main  
2 tasks: embedding *models* in unsupervised learning (*e.g.* hierarchies) and embedding  
3 *data*. To our knowledge, there are no approaches that provide embeddings for  
4 supervised models; even when hyperbolic geometry provides convenient properties  
5 for expressing popular hypothesis classes, such as decision trees (and ensembles).  
6 In this paper, we propose a full-fledged solution to the problem in three independent  
7 contributions. The first linking the theory of losses for class probability estimation  
8 to hyperbolic embeddings in Poincaré disk model. The second resolving an issue  
9 for a clean, unambiguous embedding of (ensembles of) decision trees in this model.  
10 The third showing how to smoothly tweak the Poincaré hyperbolic distance to  
11 improve its encoding and visualization properties near the border of the disk, a  
12 crucial region for our application, while keeping hyperbolicity. This last step has  
13 substantial independent interest as it is grounded in a generalization of Leibniz-  
14 Newton’s fundamental Theorem of calculus.

## 15 1 Introduction

16 Models of hyperbolic geometry have been successfully used to embed hierarchies, *i.e.* tree-based  
17 structures [14, 20, 38, 36]. Through the property of low-distortion hyperbolic embeddings [30],  
18 symbolic (hierarchies) and numeric (quality metrics) properties of unsupervised learning models can  
19 be represented in an accurate and interpretable manner [29]. When it comes to supervised learning,  
20 the current trend involves embedding *data* in a hyperbolic space, with models trained on the now  
21 hyperbolic data [9, 15, 8, 11]. It is important to note that none of these supervised methods embed  
22 *models*. Indeed, the focus of the prior literature is to learn supervised models from hyperbolic data,  
23 rather than representing supervised models in hyperbolic geometry. This is in stark contrast to the  
24 unsupervised usage described, where the (tree-based) structural properties of unsupervised models  
25 are directly exploited for good embeddings. This hints at benefits which can be used in the supervised  
26 case, especially for popular supervised models which are (structurally) tree-based.

27 **Our paper** proposes a full-fledged solution for this problem, in three independent contributions.

28 **The first** focuses on the numerical part (quality metrics) of the embedding and its link with supervised  
29 losses: to provide a natural way of embedding classification and the *confidence* of prediction, we  
30 show a link between training with the log-loss (posterior estimation) or logistic loss (real-valued  
31 classification) and hyperbolic distances computation in the Poincaré disk.

32 **The second** focuses on the symbolic part (hierarchies) of the embeddings for a popular kind of  
33 supervised of models: decision trees (and ensembles). Unlike for unsupervised learning, we show  
34 that getting an unambiguous embedding of a decision tree requires post-processing the model. Our  
35 solution extracts its *monotonic* sub-tree. This is also convenient for explainability purposes [32].

36 **The third** and most technical one focuses on visualization and the accuracy of the numerical encoding  
37 in the Poincaré disk model. The more accurate / confident is a supervised prediction, the closer to the

38 border of the disk it is represented. Thus, the best models will have their best predictions squeezed  
 39 close to the border. In addition to being suboptimal for visualization, this region is also documented  
 40 for being numerically error-prone for hyperbolic distance computation [29]. Two general fixes  
 41 currently exist: encoding values with (exponentially) more bits or utilizing a trick from Riemannian  
 42 geometry [19]. As neither is satisfactory for our usage, we propose a third, principled route: like  
 43 other distances, Poincaré distance is integral. We generalize Leibniz-Newton’s fundamental Theorem  
 44 of calculus using a generalization of classical arithmetic [22]. This generalized “tempered” integral  
 45 provides a parameter to smoothly alter the properties of the “classical” integral. When defining  
 46 the Poincarè distance, the tempering controls the hyperbolic constant of the embedding whilst also  
 47 improving the visualization and numeric accuracy of the embedded models. The generalization of  
 48 integrals to distances has independent interest as many application in ML rely on integral based  
 49 distances and distortions, *i.e.*, Bregman divergences,  $f$ -divergences, integral probability metrics, etc.  
 50 Experiments are provided on readily available domains, and all proofs, additional results and addi-  
 51 tional experiments are given in an Appendix.

## 52 2 Related work

53 Models of hyperbolic geometry have been mainly useful to embed hierarchies, *i.e.* tree-based  
 54 structures [14, 20, 38, 36], with a sustained emphasis on coding size and numerical accuracy [19, 29,  
 55 37]. In unsupervised learning and clustering, some applications have sought a simple representation  
 56 of data on the form of a tree or via hyperbolic projections [7, 6, 17, 34]. Approaches dealing with  
 57 supervised learning assumes the *data* lies in a hyperbolic space: the output visualized is primarily an  
 58 embedding of the data itself, with additional details linked to classification of secondary importance,  
 59 either support vector machines [9], neural nets, logistic regression [15], or (ensembles of) decision  
 60 trees [8, 11]. We insist on the fact that the aforementioned methods do not represent the *models* in  
 61 the hyperbolic space, even when those models tree-shaped. The question of embedding classifiers  
 62 is potentially important to improve the state of the art visualization: in the case of decision trees,  
 63 popular packages stick to a topological layer (the tree graph) to which various additional information  
 64 about classification are superimposed but without principled links to the “embedding”\*\*.

## 65 3 Definitions

66 Training a supervised model starts with a set (sample) of examples, each of which is a pair  $(\mathbf{x}, y)$ ,  
 67 where  $\mathbf{x} \in \mathcal{X}$  (a *domain*) and  $y \in \{-1, 1\}$  (labels or classes). A decision tree (DT) consists of a  
 68 binary rooted tree  $H$ , where at each internal nodes, the two outgoing arcs define a Boolean test  
 69 over observation variables (see Figure 1, center, for an example);  $\mathcal{N}(H)$  is the set of nodes of  $H$ ,  
 70  $\Lambda(H) \subseteq \mathcal{N}(H)$  is the set of leaf nodes of  $H$ .  
 71 The log-loss is best introduced in the context of the theory of *losses for class probability estimation*  
 72 (CPE) [28]. We follow the notations of [18]: A CPE loss function,  $\ell : \mathbb{Y} \times [0, 1] \rightarrow \mathbb{R}$ , is

$$\ell(y, p) \doteq [y = 1] \cdot \ell_1(p) + [y = -1] \cdot \ell_{-1}(p), \quad (1)$$

73 where  $[.]$  are Iverson brackets [16]. Functions  $\ell_1, \ell_{-1}$  are called *partial* losses. A CPE loss is  
 74 *symmetric* when  $\ell_1(p) = \ell_{-1}(1 - p), \forall p \in [0, 1]$ . The *log-loss* is a symmetric CPE loss with  
 75  $\ell_1^{\text{LOG}}(p) \doteq -\log(p)$ . The goal of learning using a CPE loss is to optimize the Bayes risk,  $\underline{L}(p) \doteq$   
 76  $\inf_{\pi} \mathbb{E}_{Y \sim p} \ell(Y, \pi)$ . In the case of the log-loss, it is the criterion used to learn DTs in C4.5 [26]:

$$\underline{L}^{\text{LOG}}(p) = -p \cdot \log p - (1 - p) \cdot \log(1 - p). \quad (2)$$

77 Models like DTs predict an empirical posterior  $\hat{\Pr}[Y = 1 | X]$  at the leaves: for an observation  $\mathbf{x}$   
 78 reaching leaf  $\lambda(\mathbf{x})$ , the posterior prediction is the local relative proportion of positive examples  
 79 in leaf  $\lambda$ , as estimated by the training sample (noted  $p_{\lambda(\mathbf{x})}^+$ ). There exists a duality between CPE  
 80 classification and the real-valued classification setting familiar to *e.g.* deep learning: optimizing  
 81 Bayes risk for the posterior is equivalent to minimizing its *convex surrogate* using the *canonical*  
 82 *link* of the loss to compute real classification [24, Theorem 1]. The canonical link of the log-loss,  
 83  $\psi_{\text{LOG}} : [0, 1] \rightarrow \mathbb{R}$  is the famed inverse sigmoid, which has a very convenient form for our purpose,

$$\psi_{\text{LOG}}(p) = \log \left( \frac{p}{1 - p} \right). \quad (3)$$

---

\*See for instance <https://github.com/parrt/dtreeviz>.

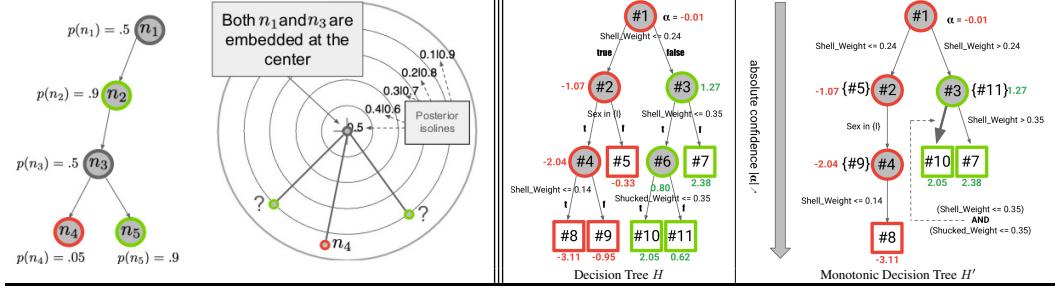


Figure 1: *Left pane:* subtree of a decision tree (DT) (colors red, green denote the majority class in each node, grey = random posterior, tests at arcs not shown,  $n_1$  is the root, no leaves shown) and its embedding following the simple recipe in (5): it is impossible to tell  $n_2$  from  $n_5$  and the tree depicted in  $\mathbb{B}_1$  is not a faithful representation of the DT nor of any of its subtrees. *Right pane:* a small DT learned on UCI abalone (left) and its corresponding monotonic decision tree (MDT, right) learned using GETMDT. In each node, the real-valued prediction ( $\psi_{\text{Log}}(p)$ ) is indicated, also in color. Observe that, indeed,  $H$  does not grant path-monotonic classification but  $H'$  does (Definition 4.1). In  $H'$ , some nodes have outdegree 1; also, internal node #6 in the DT, whose prediction is worse than its parent, disappears in  $H'$ . One arc in  $H'$  is represented with double width as its Boolean test aggregates both tests it takes to go from #3 to #10 in  $H$ . Observations that would be classified by leaf #11 (resp. #5, resp. #9) in  $H$  are now classified by internal node #3 (resp. #2, resp. #4) in  $H'$ , but predicted labels are the same for  $H$  and  $H'$  and so the accuracy of  $H$  and  $H'$  are the same.

84 Notably, the absolute value  $|\psi_{\text{Log}}(p)|$  is a *confidence* with the sign giving the predicted class. In the  
 85 case of the log-loss, the convex surrogate is hugely popular in ML: it is the logistic loss.

86 We now introduce concepts from hyperbolic geometry, particularly those from the Poincaré disk  
 87 model. Our definitions are simplified, a detailed account can be found in [5, 20]. The distance function  
 88  $d$  of a metric space is  $\tau$ -hyperbolic for some  $\tau \geq 0$  iff for any three geodesic curves  $\gamma^1, \gamma^2, \gamma^3$  linking  
 89 three points, there exists  $z$  such that  $\max_i d(z, \gamma^i) \leq \tau$ , where  $d(z, \gamma) \doteq \inf_{z' \in \gamma} d(z, z')$ . A small  
 90 hyperbolic constant guarantees thin triangles and embeddings of trees with good properties. The  
 91 Poincaré disk model,  $\mathbb{B}_1$  (negative curvature,  $-1$ ) is a popular model of hyperbolic geometry with  
 92 the hyperbolic distance between the origin and some  $z \in \mathbb{B}_1$  with Euclidean norm  $r \doteq \|z\|$  being:

$$d_{\mathbb{B}_1}(z, 0) = \log \left( \frac{1+r}{1-r} \right). \quad (4)$$

93

#### 94 4 Posterior embedding in Poincaré disk and clean embeddings for DTs

95 **Node embedding** We now exploit the similarity between log-loss confidences  $|\psi_{\text{Log}}(p)|$  and hyper-  
 96 bolic distances  $d_{\mathbb{B}_1}(z, 0)$  when embedding classification models. Notice that when given a model that  
 97 predicts posterior  $p$ , the confidence in its classification is given by the canonical link of the log-loss:

$$|\psi_{\text{Log}}(p)| = \log \left( \frac{1+r}{1-r} \right), \quad r \doteq |2p-1|. \quad (5)$$

98 If  $r = 0$ , the prediction is as bad as a fair coin's. As  $r$  edges closer to 1, prediction approaches  
 99 maximum confidence. The connection with the Poincaré distance (4) is immediate: a natural  
 100 embedding of a posterior prediction  $p$  is then a point  $z$  in the Poincaré disk  $\mathbb{B}_1$  at radius  $r \doteq \|z\|$  from  
 101 the origin (5). The origin of Poincaré disk represents the worst possible confidence. As each leaf  $\lambda$   
 102 of a DT  $H$  corresponds to a posterior  $p_{\lambda(x)}^+$ , the leaves can be embedded as described. In addition,  
 103 all nodes  $\nu \in N(H)$  in the tree also have corresponding posteriors, and thus can also be embedded  
 104 identically to leaf nodes (any node  $\nu$  in a DT was a leaf at some point during training, and hence its  
 105 corresponding leaf posterior can be used for embedding). Now, what about the arcs between nodes ?  
 106 **No clean embedding for a full DT...** We have just seen that there is a natural embedding of all nodes  
 107 of a DT in  $\mathbb{B}_1$ : embed each node with posterior  $p$  on an isoline at Euclidean distance  $r = |2p-1|$  from  
 108 the center. One could hope that naively adding the arcs of the tree could lead to clean representations  
 109 of the whole DT. However, a simple counterexample demonstrates that it cannot be the case, shown  
 110 in Figure 1 (left). In this pathological example, some distinct nodes (resp. arcs) of the DT  $H$  are  
 111 embedded in the same node (resp. edge) in  $\mathbb{B}_1$ : the depiction has no subgraph relationship to  $H$  and

112 some embedded nodes cannot be distinguish between each other without additional information.  
 113 ...but guaranteed clean embeddings of its **Monotonic part** Despite the problems of a full DT, a  
 114 “clean” embedding can be achieved by a simple post-processing of the DT. To introduce it, we define  
 115 three broad objectives for the embedding in  $\mathbb{B}_1$  of DTs and ensembles of DTs:

**Embedding objective**

- (A) The largest part of the tree in DT  $H$  is embedded in  $\mathbb{B}_1$ , it defines an injective mapping of the nodes of  $H$  and induces a subtree of  $H$  with each edge in  $\mathbb{B}_1$  corresponding to a path in  $H$ ;
- (B) *Locally*, each node  $\nu$  of  $H$  gets embedded to some  $z_\nu$  such that  $r_\nu$  is close to  $\|z_\nu\|$  (5);
- (C) *Globally*, the whole embedding remains convenient and readable to compare, in terms of confidence in classification, different subtrees in the same tree, or even between different trees. This includes their leveraging coefficient in an ensemble.

116 Our solution for clean embeddings that can satisfy (A-C) is simple in principle:

117 *Embed only the monotonic classification part of a DT,*

118 meaning for each path from the root to a leaf in  $H$ , we only embed the subsequence of nodes whose  
 119 absolute confidences are strictly increasing. To do so, we replace the DT by a path-monotonic  
 120 approximation using a new class of models we call *Monotonic Decision Trees* (MDT).

121 **Definition 4.1.** A *Monotonic Decision Tree* (MDT) is a rooted tree with a Boolean test labeling each  
 122 arc and a real-valued prediction at each node. Any sequence of nodes in a path starting from the root  
 123 is strictly monotonically increasing in absolute confidence. At any internal node, no two Boolean  
 124 tests at its outgoing arcs can be simultaneously satisfied. The classification of an observation is  
 125 obtained by the bottom-most node’s prediction reachable from the root.

126 We now introduce an algorithm which takes as input a DT  $H$  and outputs an MDT  $H'$  satisfying the  
 127 following invariants:

128 (M) For any observation  $x \in \mathcal{X}$ , the prediction  $H'(x)$  is equal to the prediction in the path followed  
 129 by  $x$  in  $H$  of its deepest node in the strictly monotonic subsequence starting from the root of  $H$ .  
 130 Figure 1 (right) presents an example of MDT  $H'$  that would be built for some DT  $H$  (left) and  
 131 satisfying (M) (unless observations have missing values,  $H'$  is unique). Figure 1 adopts some  
 132 additional conventions to ease parsing of  $H'$ :

133 (D1) Some internal nodes of  $H'$  are also tagged with labels corresponding to the leaves of  $H$ . If a  
 134 node in  $H'$  is tagged with a label of one of  $H$ ’s leaves  $\lambda$ , it indicates that examples (and predictions)  
 135 which reach  $\lambda$  in the original  $H$  are being ‘rerouted’ to  $H'$ ’s tagged node, where the original  
 136 prediction at  $\lambda$  may change in  $H'$  (some are internal nodes of  $H'$ );  
 137 (D2) Arcs in  $H'$  have a width proportional to the number of boolean tests it takes to reach its tail  
 138 from its head in  $H$ . A large width thus indicates a long path in  $H$  to improve classification confidence.

---

**Algorithm 1** GETMDT( $\nu, \mathbf{b}, \nu', \mathbb{I}$ )

**Input:** Node  $\nu \in \mathcal{N}(H)$  ( $H = \text{DT}$ ), Boolean test  $\mathbf{b}$ , Node  $\nu' \in \mathcal{N}(H')$  ( $H' = \text{MDT}$  being build from  $H$ ), forbidden posteriors  $\mathbb{I} \subset [0, 1]$ ;

```

1 : if  $\nu \in \Lambda(H)$  then
2 :   if  $p_\nu^+ \in \mathbb{I}$  then
3 :      $\nu' \leftarrow \text{TAGDTLEAF}(\nu);$                                 // tags  $\nu' \in \mathcal{N}(H')$  with info from leaf  $\nu \in \Lambda(H)$ 
4 :   else
5 :      $\nu'' \leftarrow H'.\text{NEWNODE}(\nu);$                             //  $\nu''$  will be a new leaf in  $H'$ 
6 :      $H'.\text{NEWARC}(\nu', \mathbf{b}, \nu'');$                           // adds arc  $\nu' \rightarrow_{\mathbf{b}} \nu''$  in  $H'$ 
7 :   endif
8 : else
9 :   if  $p_\nu^+ \notin \mathbb{I}$  then
10:     $\nu'' \leftarrow H'.\text{NEWNODE}(\nu);$                            //  $\nu''$  = new internal node in  $H'$ 
11:     $H'.\text{NEWARC}(\nu', \mathbf{b}, \nu'');$                           // adds arc  $\nu' \rightarrow_{\mathbf{b}} \nu''$  in  $H'$ 
12:     $\mathbb{I}_{\text{new}} \leftarrow [\min\{p_\nu^+, 1 - p_\nu^+\}, \max\{p_\nu^+, 1 - p_\nu^+\}];$  //  $\mathbb{I} \subset \mathbb{I}_{\text{new}}$ 
13:     $\nu'_{\text{new}} \leftarrow \nu'';$ 
14:     $\mathbf{b}_{\text{new}} \leftarrow \text{true};$ 
15:  else
16:     $\mathbb{I}_{\text{new}} \leftarrow \mathbb{I}; \nu'_{\text{new}} \leftarrow \nu'; \mathbf{b}_{\text{new}} \leftarrow \mathbf{b};$           //  $\nu$  yields no change in  $H'$ 
17:  endif
18:  GETMDT( $\nu.\text{leftchild}, \mathbf{b}_{\text{new}} \wedge \nu.\text{TEST}(\text{leftchild}), \nu'_{\text{new}}, \mathbb{I}_{\text{new}});$ 
19:  GETMDT( $\nu.\text{rightchild}, \mathbf{b}_{\text{new}} \wedge \nu.\text{TEST}(\text{rightchild}), \nu'_{\text{new}}, \mathbb{I}_{\text{new}});$ 
20: endif

```

---

140 It is also worth remarking that  $H'$  satisfying **(M)** has, in general, a smaller number of vertices than  $H$   
 141 but it always has the same depth. Finally, any pruning of  $H$  is a subtree of  $H$  and its corresponding  
 142 MDT is a subtree of the MDT of  $H$ . The aforementioned troubles to embed the DT in the Poincaré  
 143 disk considering **(A-C)** do not exist anymore for the MDT because the best embeddings necessarily  
 144 have all arcs going outwards in the Poincaré disk. We now present algorithm GETMDT, in Algorithm  
 145 1. To produce  $H'$ , after having initialized it to a root = single leaf, we just run

146  $\text{GETMDT}(\text{root}(H), \text{true}, \text{root}(H'), [p_{\text{root}}^+, \bar{p}_{\text{root}}^+])$

147 (we let  $p_{\text{root}}^+ \doteq \min\{p_{\text{root}}^+, 1 - p_{\text{root}}^+\}$ ,  $\bar{p}_{\text{root}}^+ \doteq \max\{p_{\text{root}}^+, 1 - p_{\text{root}}^+\}$ ). Upon finishing, the tree rooted at  
 148  $\text{root}(H')$  is the MDT sought.

149 We complete the description of GETMDT: When a leaf of  $H$  does not have sufficient confidence and  
 150 ends up being mapped to an internal node of the MDT  $H'$ , TAGDTLEAF is the procedure that tags  
 151 this internal node with information from the leaf (see **(D1)** above). We consider a tag being just the  
 152 name of the leaf (Figure 1, leaves #5, 9, 11), but other conventions can be adopted. The other two  
 153 methods we use grow MDT  $H'$  by creating a new node via NEWNODE and adding a new arc between  
 154 existing nodes via NEWARC. We easily check the following.

155 **Theorem 1.**  $H'$  built by GETMDT satisfies **(M)** with respect to DT  $H$ .

156 **Beyond DTs: embedding a sequence of trees and leveraging coefficients** The link between dis-  
 157 tances in the Poincaré model of hyperbolic geometry and the canonical link of the log-loss (5) can be  
 158 extended to leveraging coefficients in boosted combinations [31]. To get there, we need a tailored  
 159 boosting algorithm which parallels AdaBoost’s design tricks (hence different from LogitBoost [13]),  
 160 derived *in extenso* in the Appendix to save space. Assuming a basic knowledge of boosting, we give  
 161 here the main ingredients. First, the *unnormalized* weight update after adding DT  $H_j$  is

$$w_{(j+1)i} \leftarrow \frac{w_{ji}}{w_{ji} + (1 - w_{ji}) \cdot \exp(\alpha_j y_i H_j(\mathbf{x}_i))}, \quad (6)$$

162 all weights being initialized at value 1/2 at the beginning.  $\alpha_j$  is the leveraging coefficient for “weak”  
 163 classifier  $H_j$ . The final model  $\mathbf{H}_T$  is a linear combination of DTs:

$$\mathbf{H}_T(\mathbf{x}) \doteq \sum_{j=1}^T \frac{(\psi_{\text{LOG}})_j}{(\psi_{\text{LOG}})_j^*} \cdot H_j(\mathbf{x}) = \sum_{j=1}^T \frac{(\psi_{\text{LOG}})_j}{(\psi_{\text{LOG}})_j^*} \cdot \psi_{\text{LOG}}\left(p_{j\lambda(\mathbf{x})}^+\right), \quad (7)$$

164 where  $(\psi_{\text{LOG}})_j^*$  is the maximal absolute confidence of  $H_j$ :

$$(\psi_{\text{LOG}})_j^* \doteq \max_{\lambda \in \Lambda(H_j)} |\log(p_{j\lambda}^+ / (1 - p_{j\lambda}^+))| \quad (8)$$

165 (index ‘ $j$ ’ in  $p_+^+$  emphasize the use of boosting weights for posteriors) and  $p_{j\lambda(\mathbf{x})}^+$  is the posterior  
 166 estimation at leaf reached by observation  $\mathbf{x}$ .  $(\psi_{\text{LOG}})_j^*$  can be read *directly* from  $\mathbb{B}_1$  (the maximal  
 167 absolute confidence for the MDT is also that of the DT). The second leveraging part in (7) is  
 168  $(\psi_{\text{LOG}})_j \doteq \log(p_j^+ / (1 - p_j^+))$ , with  $p_j^+ \doteq (1 + r_j)/2$  and

$$r_j \doteq \mathbb{E}_{i \sim \tilde{w}_j} \left[ \frac{1}{(\psi_{\text{LOG}})_j^*} \cdot \log \left( \frac{p_{j\lambda(\mathbf{x}_i)}^+}{1 - p_{j\lambda(\mathbf{x}_i)}^+} \right) \right], \quad (9)$$

169 where  $\tilde{w}_j$  denotes weights (6) normalized. It is not hard to show that  $(\psi_{\text{LOG}})_j \geq 0$ , so it can easily be  
 170 displayed and read from  $\mathbb{B}_1$  as the confidence of an imaginary vertex with posterior  $p_j^+$  (we represent  
 171 it with a circle, see Figure 2).

172 At this stage, there only remains to detail the part which takes an MDT and embeds it in  $\mathbb{B}_1$ . Our  
 173 approach is a simple variation on Sarkar’s embedding [30] which, because of the lack of space, we  
 174 defer to the Appendix. We evaluate the quality of the MDT embedding using an expected error for  
 175 the embedding of *all* nodes:

$$\rho(H') \doteq \mathbb{E}_{\nu \sim \mathcal{N}(H')} [| |\alpha_\nu| - d_{\mathbb{B}_1}(\mathbf{0}, \mathbf{z}_\nu) | / |\alpha_\nu| ], \quad (10)$$

176 where  $\alpha_\nu$  refers to the relevant  $\psi_{\text{LOG}}(p_{j\lambda}^+)$  in (7), *i.e.* the confidence computed at the iteration  $t$  when  
 177  $\nu$  was a leaf  $\lambda$  in  $H$  (and  $\mathbf{z}_\nu \in \mathbb{B}_1$  is its embedding in  $\mathbb{B}_1$ ). Two example final representations are in  
 178 Figure 2. Notice the small errors  $\rho$  in both cases.

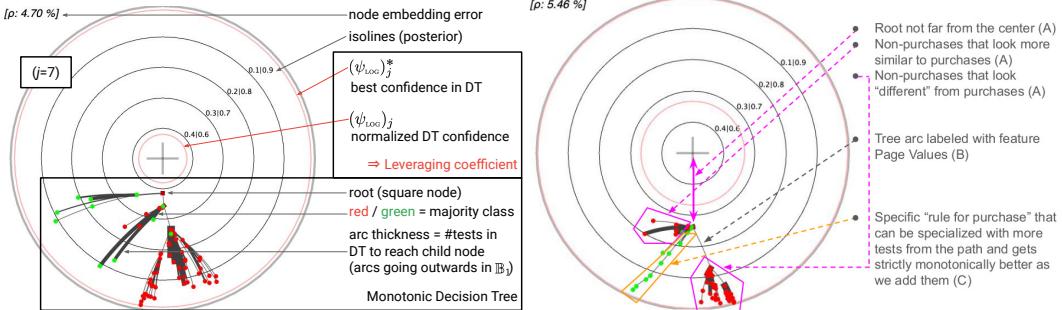


Figure 2: Two MDTs learned on UCI domain `online_shoppers_intention`. *Left:* MDT with parameter definition annotations: (10) for  $\rho$ , (8) for  $(\psi_{\log})_j^*$ , (9) for  $(\psi_{\log})_j$ . *Right:* MDT with interpretation annotations that are further detailed in Experiments, Section 6.

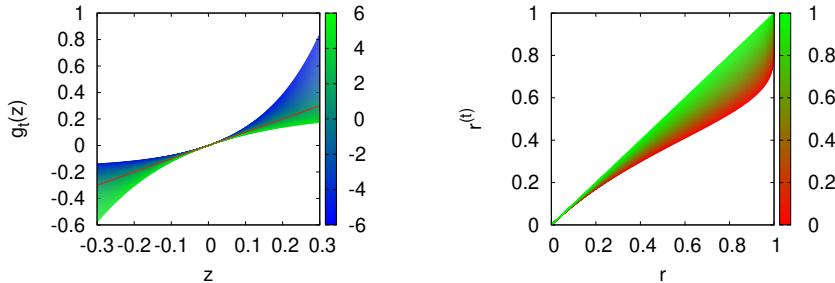


Figure 3: *Left:* plot of  $g_t(z)$  (14) for different values of  $t$  (color map on the right bar), showing where it is convex / concave. The  $t = 1$  case ( $g_1(z) = z$ ) is emphasized in red. *Right:* suppose  $r \doteq \|z\|$  is the Euclidean norm of a point  $z$  in Poincaré disk  $\mathbb{B}_1$ . Fix  $t \in [0, 1]$  (color map on the right bar). The plot gives the (Euclidean) norm  $r^{(t)}$  of a point  $z^{(t)}$  in the  $t$ -self such that  $d_{\mathbb{B}_1}(z^{(t)}, \mathbf{0}) = d_{\mathbb{B}_1}(z, \mathbf{0})$  (18). Remark that even for  $r$  very close to 1 we can have  $r^{(t)}$  substantially smaller (e.g.  $r^{(t)} < 0.8$ ).

## 179 5 Smoothly altering integrals: T-calculus and the t-self of Poincaré disk

180 It is apparent from Figure 2 (right, lowest red pentagon), that when utilizing standard hyperbolic  
181 distances  $d_{\mathbb{B}_1}(z, \mathbf{0})$  (4), the best MDT nodes are embedded close to the border. In addition to  
182 obviously not being great for visualization, these nodes are at risk of having numeric error “push” the  
183 nodes to the border  $\partial\mathbb{B}_1$ , thereby giving a false depiction of infinite confidence. In this section, we  
184 provide alternative distances to prevent this phenomena. First, we quantify the numerical risk, which  
185 has been defined by the critical region where high numeric error can occur [19, 29].

186 **Definition 5.1.** A point  $x \in \mathbb{B}_1$  is said  $k$ -close to boundary  $\partial\mathbb{B}_1$  if  $\|x\| = 1 - 10^{-k}$ . It is said  
187 encodable iff  $\|x\| < 1$  in machine encoding (it is not “pushed to the boundary”).

188 Machine encoding constrains the maximum possible  $k$ : in the double-precision floating-point rep-  
189 resentation (Float64),  $k \approx 16$  [19]. In the case of Poincaré disk, the maximal distance  $d_*$  from the  
190 origin  $\mathbf{0}$  that this authorizes (before numerical error “warps” a point to the border) is a small affine  
191 order in  $k$  [19]:

$$d_* \leq \log(2) + \log(10) \cdot k + O(10^{-k}), \quad (11)$$

192 Hence, in practice, only a ball of radius  $d^* \approx 38$  around the origin can be accurately represented. This  
193 is in deep contrast with Euclidean representation, where  $d_* = \Omega(2^k)$ . We now provide a principled  
194 solution to changing  $d^*$  while keeping hyperbolicity, which relies on a crucial generalization of  
195 Leibniz-Newton’s fundamental Theorem of calculus.

196 **T-calculus** We let  $[n] \doteq \{1, 2, \dots, n\}$  for  $n \in \mathbb{N}_{>0}$ . At the core of our generalization is the replacement  
197 of the addition by the tempered (or t)-addition [22],  $z \oplus_t z' \doteq z + z' + (1-t)zz'$ . The additional term  
198 is a scaled saddle point curve  $zz'$ , positive when the sign of  $z$  and  $z'$  agree and negative otherwise.  
199 Hereafter,  $f$  is defined on an interval  $[a, b]$ . We define a generalization of Riemann integration  
200 to t-algebra (see [2] for a preliminary development which does not provide the generalization of  
201 Leibniz-Newton’s fundamental Theorem of calculus) and for this objective, given an interval  $[a, b]$   
202 and a division  $\Delta$  of this interval using  $n + 1$  reals  $x_0 \doteq a < x_1 < \dots < x_{n-1} < x_n \doteq b$ , we define

203 the Riemann  $t$ -sum of  $f$  over  $[a, b]$  using  $\Delta$ , for a set  $\{\xi_i \in [x_{i-1}, x_i]\}_{i \in [n]}$ ,

$$S_{\Delta}^{(t)}(f) \doteq (\bigoplus_{i \in [n]} (x_i - x_{i-1}) \cdot f(\xi_i)) \quad (12)$$

204 ( $S_{\Delta}^{(1)}(f)$  is the classical Riemann summation). Let  $s(\Delta) \doteq \max_i |\mathbb{I}_i|$  denote the step of division  $\Delta$ .  
205 The conditions for  $t$ -Riemann integration are the same as for  $t = 1$ .

206 **Definition 5.2.** Fix  $t \in \mathbb{R}$ . A function  $f$  is  $t$ -(Riemann) integrable over  $[a, b]$  iff  $\exists L \in \mathbb{R}$  such that  
207  $\forall \varepsilon > 0, \exists \delta > 0, \forall$  division  $\Delta$  with  $s(\Delta) < \delta, |S_{\Delta}^{(t)}(f) - L| < \varepsilon$ . When this happens, we note

$$\int_a^b f(x) d_t x = L. \quad (13)$$

208 The case  $t = 1$ , Riemann integration, is denoted using classical notations. We now prove our first  
209 main results for this Section, namely the link between  $t$ -Riemann integration and Riemann integration.

210 **Theorem 2.** Any function is either  $t$ -Riemann integrable for all  $t \in \mathbb{R}$  simultaneously, or for none. In  
212 the former case, we have the relationship:  $(\log_t \doteq (z^{1-t} - 1)/(1-t)$  for  $t \neq 1$  and  $\log_1 \doteq \log$

$$\int_a^b f(u) d_t u = \mathfrak{g}_t \left( \int_a^b f(u) du \right), \forall t \in \mathbb{R}, \quad \text{with } \mathfrak{g}_t(z) \doteq \log_t \exp z. \quad (14)$$

213 For  $t = 0$ , we get  $1 + {}^{(0)}\int_a^b f(u) d_t u = \exp \int_a^b f(u) du$ , which happens to be Volterra's integral [33,  
214 Theorem 5.5.11]. Classical Riemann integration and derivation are fundamental inverse operations.  
215 The classical derivation is sometimes called "Euclidean derivative" [4]. We now elicit the notion of  
216 derivative, which is the "inverse" of  $t$ -Riemann integration. Unsurprisingly, it generalizes Euclidean  
217 derivative. The Theorem stands as a generalization of the classical fundamental Theorem of calculus.

218 **Theorem 3.** Let  $z \ominus_t z' \doteq (z - z')/(1 + (1-t)z')$  denote the tempered subtraction. When it exists,  
219 define  $D_t f(z) \doteq \lim_{\delta \rightarrow 0} (f(z + \delta) \ominus_t f(z)) / \delta$ . Suppose  $f$  to be  $t$ -Riemann integrable. Then, the  
220 function  $F$  defined by  $F(z) \doteq {}^{(t)}\int_a^z f(u) d_t u$  is such that  $D_t F = f$ . We call  $F$  a  $t$ -primitive of  $f$   
221 (which zeroes when  $z = a$ ) and  $D_t F$  the  $t$ -derivative of  $F$ .

223 The function  $\mathfrak{g}_t$  (Figure 3) is key to many of our results; it has quite remarkable properties: in particular,  
224 it is strictly increasing for any  $t \in \mathbb{R}$ , strictly concave for any  $t > 1$ , strictly convex for any  $t < 1$  and  
225 such that  $\text{sign}(\mathfrak{g}_t(z)) = \text{sign}(z), \forall t \in \mathbb{R}$ . One can also note that  $D_t \mathfrak{g}_t(z) = 1$ . The Appendix proves  
226 these results and provides many more, showing how Theorems 2 and 3 naturally "percolate" through  
227 many properties known for classical integration and the fundamental Theorem of calculus.

228 **The  $t$ -self of Poincaré disk** Let us now put T-calculus to good use in our context. For a set  $\mathcal{X}$   
229 endowed with function  $d : \mathcal{X} \times \mathcal{X} \rightarrow [0, +\infty]$  (e.g. Poincaré disk and the Poincaré hyperbolic  
230 distance), the  $t$ -self of  $\mathcal{X}$  is the set (implicitly) endowed with  $d^{(t)} \doteq \mathfrak{g}_t \circ d$ . Before going further, let  
231 us provide the elegant  $d^{(t)}$  associated to Poincaré disk's  $t$ -self (with  $r \doteq \|z\|$  as in (4)):

$$d_{\mathbb{B}_1}^{(t)}(\mathbf{0}, z) = \log_t \exp d_{\mathbb{B}_1}(\mathbf{0}, z) = \log_t \left( \frac{1+r}{1-r} \right). \quad (15)$$

232 The following Lemma (shown in Appendix) demonstrates the usefulness of T-calculus to address the  
233 encoding problem in  $\mathbb{B}_1$ .

234 **Lemma 1.** In Poincaré disk model, pick any increasing function  $g(k) \geq 0$ . For the choice  $t =$   
235  $1 - f(k)$  where  $f(k) \in \mathbb{R}$  is any function satisfying

$$\frac{\log(1 + f(k)g(k))}{f(k)} \leq \log(10) \cdot k, \quad (16)$$

236 we get in lieu of (11) the maximal  $d_*^{(t)}$  of  $d^{(t)}$  satisfying  $d_*^{(t)} \geq g(k)$ , and the new hyperbolic constant  
237  $\tau_t$  of the  $t$ -self satisfies

$$\tau_t = \frac{\exp(f(k)\tau) - 1}{f(k)}. \quad (17)$$

238 Since both  $\log(1 + x) = x + o(x)$  and  $\exp(x) - 1 = x + o(x)$  in a neighborhood of 0, we  
239 check that (16) and (17) get back to properties of the Poincaré model as  $f(k) \rightarrow 0$  [19]. We then

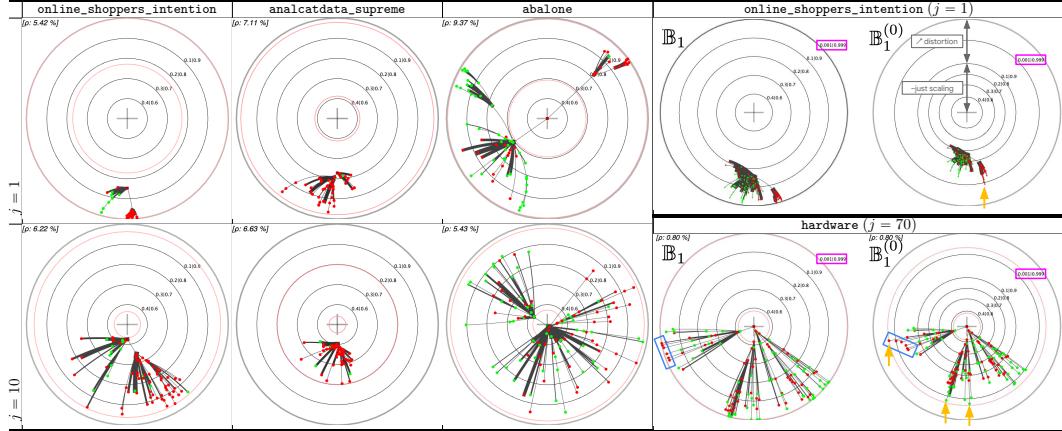


Table 1: *Left pane*: embedding in Poincaré disk of the MDTs corresponding to the DTs learned at the 1<sup>st</sup> (top row) and 10<sup>th</sup> (bottom row) boosting iteration, on four UCI domains. Stark differences emerge between these domains from the plots alone. *Right pane*: comparison between Poincaré disk embedding and its  $t = 0$  t-self for an MDT learned on UCI `online_shoppers_intention` (top, boosting coefficients information not shown) and UCI `hardware` (bottom). The  $p \in \{0.001, 0.999\}$  isoline (`rectangle`) is barely distinguishable from  $\partial \mathbb{B}_1$  but is clearly distinct from  $\partial \mathbb{B}_1^{(0)}$ . Note that in  $\mathbb{B}_1^{(0)}$ , the center looks similar to a scaling (zoom) of  $\mathbb{B}_1$ ; while near the border, the high nonlinearity of  $\mathbb{B}_1^{(0)}$  allows us to spot nodes that have high confidence / training accuracy (in `orange`) but can hardly be distinguished from the bulk of “just good” nodes in  $\mathbb{B}_1$ . Note also the set of `red` nodes in the Poincaré disk for  $j = 70$  (`rectangle`) that mistakenly look aligned, but not in the t-self.

240 have two cases: we can easily approach back the Euclidean  $d_*$  by picking  $f(k) > 0$ : choosing  
241  $f(k) = 1$  gets there and we still keep a finite hyperbolic constant, albeit exponential in the former  
242 one. If however we want to improve further the hyperbolic constant, we can pick some admissible  
243  $f(k) < 0$ , but then (16) will constrain  $g(k)$  to a very small value. Note that depending on the sign  
244 of  $f(k)$ , the triangle inequality can still hold or be replaced by a weaker version since we can show  
245  $d_{\mathbb{B}_1}^{(t)}(\mathbf{x}, \mathbf{z}) \leq d_{\mathbb{B}_1}^{(t)}(\mathbf{x}, \mathbf{y}) + d_{\mathbb{B}_1}^{(t)}(\mathbf{y}, \mathbf{z}) + \max\{0, f(k)\} \cdot d_{\mathbb{B}_1}^{(t)}(\mathbf{x}, \mathbf{y}) \cdot d_{\mathbb{B}_1}^{(t)}(\mathbf{y}, \mathbf{z})$  in  $\mathbb{B}_1^{(t)}$ .

246 Let us finally investigate how we can perform a general embedding in the  $t$ -self  $\mathbb{B}_1^{(t)}$  of the Poincaré  
247 disk to be (i) more encoding-savvy and (ii) have points originally close to  $\partial \mathbb{B}_1$  substantially further  
248 away from the  $t$ -self  $\mathbb{B}_1^{(t)}$  border. We also add a third constraint (iii) that the distortions must be *fair*  
249 w.r.t. the original Poincaré disk model embedding; *i.e.* if some  $\mathbf{z} \in \mathbb{B}_1$  (*e.g.* the coordinate of a  
250 DT node) gets mapped to  $\mathbf{z}^{(t)} \in \mathbb{B}_1^{(t)}$  (the t-self), then we request

$$d_{\mathbb{B}_1}^{(t)}(\mathbf{z}^{(t)}, \mathbf{0}) = d_{\mathbb{B}_1}(\mathbf{z}, \mathbf{0}). \quad (18)$$

251 All three conditions point to a non-linear embedding  $\varphi_t : \mathbb{B}_1 \rightarrow \mathbb{B}_1^{(t)}$ . The  $t$ -self offers a simple  
252 convenient solution with a trivial design for  $\varphi_t$ : We compute  $\mathbf{z}^{(t)}$  by a simple scaling of the Euclidean  
253 norm of  $\mathbf{z}$  in  $\mathbb{B}_1$  to ensure that (18) holds. Figure 3 (right) displays the corresponding relationship  
254 between  $r$  and  $r^{(t)}$  when  $t \in [0, 1]$ , clearly achieving (ii) in addition to (iii). For (i), even for  $t < 1$   
255 very close to 1 and  $\mathbf{z}$  close to  $\partial \mathbb{B}_1$  ( $r$  close to 1), the mapping can send  $\mathbf{z}^{(t)}$  substantially “back in”  
256 the t-self  $\mathbb{B}_1^{(t)}$ : for  $t = 0.7$  and  $r = 1 - 10^{-4}$  – *i.e.*  $k = 4$  in Definition 5.1 – we get  $r^{(t)} \approx 0.96$ ,  
257 authorizing encoding with a less “risky”  $k = 2$ . One additional benefit from (18), visible from Figure  
258 3 (right) is that the distortion is low near the center of the disk.

## 259 6 Experiments

260 We summarize a number of experiments (Table 1), provided otherwise *in extenso* in the Appendix. In  
261 the top-down induction scheme for DT, the leaf chosen to be split is the heaviest non pure leaf, *i.e.*  
262 the one among the leaves containing both classes with the largest proportion of training examples.  
263 Induction is stopped when the tree reaches a fixed size, or when all leaves are pure, or when no further  
264 split allows decreasing the expected log-loss. We do not prune DTs. Our domains are public (*Cf*

265 Appendix).

266 **Poincaré embeddings of DTs / MDTs** See Figure 2 (left) for a summary of the visualization. We  
267 remind that the center of the disk is “poor” classification (confidence 0), or, in the context of boosting,  
268 random classification. A striking observation is the sheer variety of patterns that are plainly obvious  
269 from our visualization but would otherwise be difficult to spot using classical tools. Some are  
270 common to all domains: as boosting iterations increase, low-depth tree nodes tend to converge to the  
271 center of the disk, indicating increased hardness in classification. This is the impact of a well known  
272 effect of boosting, whose weight modifications make the problem harder. Among domain-dependent  
273 patterns, highly imbalanced domains get a root predictably initially embedded “far” from the origin  
274 (`online_shoppers_intention`, `analcatdata_supreme`) while more balanced domains get their  
275 roots embedded near the origin (`abalone`). Across the experiment, all trees have at most 200 nodes:  
276 while this clearly provides very good models after a large number of iterations on some domains,  
277 it is not enough to get good models after just a few iterations on others (`analcatdata_supreme`).  
278 Interpreting the hyperbolic embeddings can also be telling: consider Figure 2 (right), which a MDT of  
279 the ( $j = 3$ ) boosted DT learned on `online_shoppers_intention`. Notice the low embedding error  
280 (5.46%). We see in (A) (magenta) that classes are more balanced at the root compared to previous  
281 trees (All DTs in Appendix, Table III) because the root is closer to the center of the disk: thus, hard  
282 examples belong to both classes (no class becomes much harder to “learn” than the other). Two  
283 clearly distinct subtrees are associated to non-purchase patterns (red, A). The bottom-most subtree  
284 achieves very good confidence with several leaves close to the border: these are non-purchase patterns  
285 clearly distinct from purchase patterns (green nodes); the whole subtree dangles from the root via  
286 a test on feature `PageValues` (grey, B) achieving a substantial boost in confidence; many nodes are  
287 way past posterior isoline  $p \in \{0.1, 0.9\}$ . Red subtree in (A) is a lot closer to purchase patterns (C).  
288 One distinct pattern of purchase emerges, built from tests that always *strictly* increase its prediction’s  
289 confidence (orange). The full rule achieves very high accuracy (leaf posterior nears 0.99).

290 **Interest of the t-self for visualization** As demonstrated by Table 1 (right) and Appendix, the t-self  
291 is particularly useful to tell the best parts of a tree. Because it also manages to “push back” the bulk  
292 of nodes from the border  $\partial\mathbb{B}_1^{(t)}$ , it displays the t-self might also be useful as a *standard* encoding (use  
293 (18) to access Poincaré disk quantities and embedding).

294 **MDT vs DT classification** We traded the hyperbolic visualization of a DT – bound to be poor – for  
295 the visualization of its “monotonic part”, an MDT. MDT visualization provides lots of clues about the  
296 DT as well, but an interesting question is how deep this relationship runs: can an MDT be equivalent  
297 to its DT, classification-wise ? This question deserves a formal, extensive answer of its own. However,  
298 we ran a simple experiment (Appendix) displaying that, in the context of our experiments, similar  
299 predictions happens to be statistically true for a large majority of our domains.

## 300 7 Conclusion

301 This paper proposes three separate contributions that altogether provide a solution for hyperbolic  
302 embedding of (ensembles of) decision trees (DT), via (i) a link between losses for class probability  
303 estimation and hyperbolic distances, (ii) the design of a clean, unambiguous embedding of a DT  
304 via its monotonic subtree and (iii) a way to improve visualization and encoding properties of the  
305 Poincaré disk model using its t-self, a notion we introduce. Each of these contributions are indeed  
306 separate: for example, one could reuse (i) to embed models different from DTs or (iii) to perform  
307 hyperbolic embeddings of objects being not even related to supervised models. From this standpoint,  
308 we believe that the way we address (iii) opens general applications even outside hyperbolic geometry.  
309 Indeed, we generalize classical integration and derivation to a context encompassing the concept of  
310 additivity, upon which integration is built, extending standard properties of integration and derivation  
311 in a natural way (see the Appendix for many examples). That such properties can be obtained without  
312 additional technical “effort” bodes well for perspectives of developments and applications in other  
313 subfields of ML, where such tools could be used, not just in our geometric context, to smoothly tweak  
314 distortion measures. Many distortion measures used in ML are indeed integrals in nature: Bregman  
315 divergences,  $f$ -divergences, integral probability metrics, etc.

316 One inherent limitation of our work is linked to the use of DTs: they typically fit very well to tabular  
317 data (attribute-value data) but otherwise should be used with caution. Our work is a first step towards  
318 more sophisticated visualization for supervised models (Section 2), looking for tailored fit geometric  
319 spaces to best blend their symbolic and numerical properties. More can be done and more has to be  
320 done: at an age of data collection and ML compute still ramping up, seeking model “pictures” worth  
321 a thousand words is a challenge but a necessity for responsible AI and explainability.

322 **References**

- 323 [1] E. Amid, F. Nielsen, R. Nock, and M.-K. Warmuth. Optimal transport with tempered exponential  
324 measures. In *AAAI'24*, 2024.
- 325 [2] Ehsan Amid, Frank Nielsen, Richard Nock, and Manfred K Warmuth. The tempered Hilbert  
326 simplex distance and its application to non-linear embeddings of TEMs. *arXiv preprint arXiv:2311.13459*, 2023.
- 328 [3] A. Banerjee, S. Merugu, I. Dhillon, and J. Ghosh. Clustering with Bregman divergences. In  
329 *Proc. of the 4<sup>th</sup> SIAM International Conference on Data Mining*, pages 234–245, 2004.
- 330 [4] A.-F. Beardon and D. Minda. The hyperbolic metric and geometric function theory. In *Proc. of  
331 the International Workshop on Quasiconformal Mappings and their Applications*, 2005.
- 332 [5] B.-H. Bowditch. *A course on geometric group theory*. Math. Soc. Japan Memoirs, 2006.
- 333 [6] I. Chami, A. Gu, V. Chatziafratis, and C. Ré. From trees to continuous embeddings and back:  
334 Hyperbolic hierarchical clustering. In *NeurIPS\*33*, 2020.
- 335 [7] I. Chami, A. Gu, D. Nguyen, and C. Ré. Horopca: Hyperbolic dimensionality reduction via  
336 horospherical projections. In *38<sup>th</sup> ICML*, volume 139 of *Proceedings of Machine Learning  
337 Research*, pages 1419–1429. PMLR, 2021.
- 338 [8] P. Chlenski, E. Turok, A. Khalil Moretti, and I. Pe'er. Fast hyperboloid decision tree algorithms.  
339 In *12<sup>th</sup> ICLR*, 2024.
- 340 [9] H. Cho, B. Demeo, J. Peng, and B. Berger. Large-margin classification in hyperbolic space. In  
341 *22<sup>nd</sup> AISTATS*, volume 89 of *Proceedings of Machine Learning Research*, pages 1832–1840.  
342 PMLR, 2019.
- 343 [10] Z. Cranko and R. Nock. Boosted density estimation remastered. In *36<sup>th</sup> ICML*, pages 1416–  
344 1425, 2019.
- 345 [11] L. Doorenbos, P. Márquez-Neila, R. Sznitman, and P. Mettes. Hyperbolic random forests. *CoRR*,  
346 abs/2308.13279, 2023.
- 347 [12] D. Dua and C. Graff. UCI machine learning repository, 2021.
- 348 [13] J. Friedman, T. Hastie, and R. Tibshirani. Additive Logistic Regression : a Statistical View of  
349 Boosting. *Ann. of Stat.*, 28:337–374, 2000.
- 350 [14] O.-E. Ganea, G. Bécigneul, and T. Hofmann. Hyperbolic entailment cones for learning hierar-  
351 chical embeddings. In *35<sup>th</sup> ICML*, volume 80 of *Proceedings of Machine Learning Research*,  
352 pages 1632–1641. PMLR, 2018.
- 353 [15] O.-E. Ganea, G. Bécigneul, and T. Hofmann. Hyperbolic neural networks. In *NeurIPS\*31*,  
354 pages 5350–5360, 2018.
- 355 [16] D.-E. Knuth. Two notes on notation. *The American Mathematical Monthly*, 99(5):403–422,  
356 1992.
- 357 [17] M.-T. Law, R. Liao, J. Snell, and R.-S. Zemel. Lorentzian distance learning for hyperbolic  
358 representations. In *36<sup>th</sup> ICML*, volume 97 of *Proceedings of Machine Learning Research*, pages  
359 3672–3681. PMLR, 2019.
- 360 [18] Y. Mansour, R. Nock, and R.-C. Williamson. Random classification noise does not defeat all  
361 convex potential boosters irrespective of model choice. In *40<sup>th</sup> ICML*, 2023.
- 362 [19] G. Mishne, Z. Wan, Y. Wang, and S. Yang. The numerical stability of hyperbolic representation  
363 learning. In *40<sup>th</sup> ICML*, 2023.
- 364 [20] M. Nickel and D. Kiela. Poincaré embeddings for learning hierarchical representations. In  
365 *NeurIPS\*30*, pages 6338–6347, 2017.
- 366 [21] F. Nielsen and R. Nock. On Rényi and Tsallis entropies and divergences for exponential families.  
367 *CoRR*, abs/1105.3259, 2011.
- 368 [22] L. Nivanen, A. Le Méhauté, and Q.-A. Wang. Generalized algebra within a nonextensive  
369 statistics. *Reports on Mathematical Physics*, 52:437–444, 2003.
- 370 [23] R. Nock, W. Bel Haj Ali, R. D'Ambrosio, F. Nielsen, and M. Barlaud. Gentle nearest neighbors  
371 boosting over proper scoring rules. *IEEE Trans.PAMI*, 37(1):80–93, 2015.
- 372 [24] R. Nock and A. K. Menon. Supervised learning: No loss no cry. In *37<sup>th</sup> ICML*, 2020.

- 373 [25] M.-C. Pardo and I. Vajda. About distances of discrete distributions satisfying the data processing  
 374 Theorem of Information Theory. *IEEE Trans. IT*, 43:1288–1293, 1997.
- 375 [26] J. R. Quinlan. *C4.5 : programs for machine learning*. Morgan Kaufmann, 1993.
- 376 [27] J.-G. Ratcliffe. *Foundations of Hyperbolic Manifolds*. Springer Graduate Texts in Mathematics,  
 377 1994.
- 378 [28] M.-D. Reid and R.-C. Williamson. Information, divergence and risk for binary experiments.  
 379 *JMLR*, 12:731–817, 2011.
- 380 [29] F. Sala, C. De Sa, A. Gu, and C.Ré. Representation tradeoffs for hyperbolic embeddings. In *35<sup>th</sup>  
 381 ICML*, volume 80 of *Proceedings of Machine Learning Research*, pages 4457–4466. PMLR,  
 382 2018.
- 383 [30] R. Sarkar. Low distortion Delaunay embedding of trees in hyperbolic plane. In *GD'11*, volume  
 384 7034, pages 355–366. Springer, 2011.
- 385 [31] R. E. Schapire and Y. Singer. Improved boosting algorithms using confidence-rated predictions.  
 386 *MLJ*, 37:297–336, 1999.
- 387 [32] Andrew D Selbst and Solon Barocas. The intuitive appeal of explainable machines. *Fordham L.  
 388 Rev.*, 87:1085, 2018.
- 389 [33] A. Slavik. *Product integration, its history and applications*. Matfyzpress, Praze, 2007.
- 390 [34] R. Sonthalia and A.-C. Gilbert. Tree! I am no tree! I am a low dimensional hyperbolic  
 391 embedding. In *NeurIPS\*33*, 2020.
- 392 [35] T. van Erven and P. Harremoës. Rényi divergence and kullback-leibler divergence. *IEEE Trans.  
 393 IT*, 60:3797–3820, 2014.
- 394 [36] M. Yang, M. Zhou, R. Ying, Y. Chen, and I. King. Hyperbolic representation learning: Revisiting  
 395 and advancing. In *40<sup>th</sup> ICML*, volume 202 of *Proceedings of Machine Learning Research*,  
 396 pages 39639–39659, 2023.
- 397 [37] T. Yu and C. De Sa. Numerically accurate hyperbolic embeddings using tiling-based models.  
 398 In *NeurIPS\*32*, pages 2021–2031, 2019.
- 399 [38] T. Yu, T.-J.-B. Liu, A. Tseng, and C. De Sa. Shadow cones: Unveiling partial orders in  
 400 hyperbolic space. *CoRR*, abs/2305.15215, 2023.

# Supplementary Material

## Abstract

402 This is the Supplementary Material to Paper "Hyperbolic Embeddings of Super-  
403 vised Models" submitted to NeurIPS'24.

404 To differentiate with the numberings in the main file, the numbering of Theorems, etc. is letter-based  
405 (A, B, ...).

406 **Table of contents**

407 <b>Broader Impact</b>	Pg 13
408	
409 <b><i>t</i>-algebra and <i>t</i>-additivity</b>	Pg 13
410	
411 <b>Supplementary material on proofs</b>	Pg 14
412	
413     → Proof of Theorem 2	Pg 14
414     → Proof of Theorem 3	Pg 16
415     → Additional and helper results for Theorems 2 and 3	Pg 16
416     → Sets endowed with a distortion and their <i>t</i> -self: statistical information	Pg 19
417     → Using the T-calculus in the Lorentz model of hyperbolic geometry	Pg 20
418     → Proof of Lemma 1	Pg 21
419     → Proof of Theorem 1	Pg 21
420     → Boosting with the logistic loss <i>à-la</i> AdaBoost	Pg 22
421     → Modifying Sarkar's embedding in Poincaré disk	Pg 23
422	
423 <b>Supplementary material on experiments</b>	Pg 23
424	
425     → Domains	Pg 23
426     → Visualizing a DT <i>via</i> its MDT	Pg 23
427     → All Poincaré disk embeddings	Pg 25
428     → Interest of the <i>t</i> -self for visualization	Pg 25
429	

430 **A Impact Statement**

431 This paper presents work whose goal is to advance the field of Machine Learning. Beyond our  
 432 proposed tempered calculus, the monotonic decision trees and  $t$ -self hyperbolic embedding contribute  
 433 to the field of visualization and explainable AI for ML models. From a societal perspective, these  
 434 visualizations will provide better methods for deployed decision tree models to be scrutinized. A  
 435 potential negative of our approach is the required reduction from normal decision trees to monotonic  
 436 decision trees, where the visualization does not directly correspond to the initial decision tree. As  
 437 such, practitioners should be careful when extracting inferences from the reduced monotonic decision  
 438 tree. Nevertheless, we show that – at least empirically within the bounds of our settings – that  
 439 monotonic decision trees are similar in prediction to their original decision tree counterparts.

440 **B  $t$ -algebra and  $t$ -additivity**

441 We provide here a few more details on the basis of our paper, the  $t$ -algebra and the  $t$ -additivity of  
 442 some divergences.

443 **B.I Primer on  $t$ -algebra**

444 Classical arithmetic over the reals can be used to display duality relationships between operators using  
 445 the log, exp functions, such as for example  $\log(a/b) = \log a - \log b$ ,  $\exp(a+b) = \exp(a) \cdot \exp(b)$ ,  
 446 and so on. They can also be used to define one operator from another one. There is no difference  
 447 between the operators appearing inside and outside functions. In the  $t$ -algebra, a difference appears  
 448 and such relationships can be used to define the  $t$ -operators from those over the reals, as indeed one  
 449 can define the tempered addition

$$x \oplus_t y \doteq \log_t(\exp_t(x) \cdot \exp_t(y)),$$

450 the tempered subtraction,

$$x \ominus_t y \doteq \log_t(\exp_t(x)/\exp_t(y)),$$

451 (both simplifying to the expressions we use), and of course the tempered product and division,

$$x \otimes_t y \doteq \exp_t(\log_t(x) + \log_t(y)) \quad ; \quad x \oslash_t y \doteq \exp_t(\log_t(x) - \log_t(y)),$$

452 whose simplified expression appears e.g. in [1]. See also **(author?)** [22].

453 **B.II Functional  $t$ -additivity**

454 As is well-known, Boltzman-Gibbs (and so, Shannon) entropy is additive while Tsallis entropy is  
 455  $t$ -additive. Note that additivity for BG requires being on the simplex, but  $t$ -additivity of Tsallis  
 456 technically requires only positive measures – the simplex restriction ensures the limit exists for  $t \rightarrow 1$   
 457 and then T→BG.

458 **Divergences can also be  $t$ -additive** The KL divergence

$$D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q}) \doteq \sum_k p_k \log(p_k/q_k)$$

459 is additive on the simplex but not  $t$ -additive: using a decomposition of  $\mathbf{p}, \mathbf{q}$  as product of (independent)  
 460 distributions ( $\mathbf{p}_1, \mathbf{p}_2$  and  $\mathbf{q}_1, \mathbf{q}_2$ ) using the cartesian product of their support, we indeed check  
 461  $D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q}) = D_{\text{KL}}(\mathbf{p}_1 \parallel \mathbf{q}_1) + D_{\text{KL}}(\mathbf{p}_2 \parallel \mathbf{q}_2)$  with  $p_{ij} = p_{1i}p_{2j}$  and  $q_{ij} = q_{1i}q_{2j}$ . On the other hand,  
 462 Tsallis divergence [21] is  $t$ -additive on positive measures with such a decomposition, with

$$D_T(\mathbf{p} \parallel \mathbf{q}) \doteq \frac{\sum_k p_k (p_k/q_k)^{1-t} - 1}{1-t},$$

463 and we check that  $D_T(\mathbf{p} \parallel \mathbf{q}) = D_T(\mathbf{p}_1 \parallel \mathbf{q}_1) + D_T(\mathbf{p}_2 \parallel \mathbf{q}_2) + (1-t) \cdot D_T(\mathbf{p}_1 \parallel \mathbf{q}_1) \cdot D_T(\mathbf{p}_2 \parallel \mathbf{q}_2)$   
 464 (additional requirement to be on the simplex for convergence to  $D_{\text{KL}}$  as  $t \rightarrow 1$ ). For  $t \notin (1, 2)$  and on  
 465 the simplex, Tsallis divergence is an  $f$ -divergence with generator  $z \mapsto (z^{2-t} - 1)/(1-t)$ . Similarly,  
 466 the tempered relative entropy is  $t$ -additive on the co-simplex, where in this case

$$D_t(\mathbf{p} \parallel \mathbf{q}) \doteq \frac{1 - \sum_k p_k q_k^{1-t}}{1-t}, \mathbf{p}, \mathbf{q} \in \tilde{\Delta}_m.$$

467 The tempered relative entropy is a Bregman divergence with generator  $z \mapsto z \log_t z - \log_{t-1} z$ .

468 **C Supplementary material on proofs**

469 **C.I Proof of Theorem 2**

470 The Theorem is tautological for  $t = 1$  so we prove it for  $t \neq 1$ . Denote  $\mathcal{P}(\mathcal{S})$  the set of subsets of set  
 471  $\mathcal{S}$  and  $\mathcal{P}_*(\mathcal{S}) \doteq \mathcal{P}(\mathcal{S}) \setminus \{\emptyset\}$ . We first transform  $S_{\Delta_n}^{(t)}(f)$  (index  $n$  shown for readability) in a better  
 472 suited expression:

$$\begin{aligned} S_{\Delta_n}^{(t)}(f) &\doteq (\bigoplus_t)_{i=1}^n |\mathbb{I}_i| f(\xi_i) \\ &= \sum_{P \in \mathcal{P}_*([n])} (1-t)^{|P|-1} \cdot \prod_{i \in P} |\mathbb{I}_i| f(\xi_i) \\ &= \frac{1}{1-t} \cdot \sum_{P \in \mathcal{P}_*([n])} \prod_{i \in P} (1-t) |\mathbb{I}_i| f(\xi_i) \\ &= \frac{1}{1-t} \cdot \left( \prod_{i=1}^n (1 + (1-t) |\mathbb{I}_i| f(\xi_i)) - 1 \right), \end{aligned} \quad (19)$$

473 where we have used  $\mathbb{I}_i \doteq [x_{i-1}, x_i]$  for conciseness. We now need a technical Lemma.

474 **Lemma 2.** Fix any  $0 \leq v < 8/10$  and consider any  $n$  reals  $q_i, i \in [n]$  such that  $|q_i| \leq v, \forall i \in [n]$ .  
 475 Then

$$1 \leq \frac{(1 + \frac{1}{n} \cdot \sum_i q_i)^n}{\prod_i (1 + q_i)} \leq \exp(nv \cdot v). \quad (20)$$

476 *Proof.* Suppose all the  $q_i, i \in [n]$  satisfy  $q_i \in [u, v]$  and let  $\varphi$  be strictly convex differentiable and  
 477 defined over  $[u, v]$ . Then it comes from (**author?**) [10, Lemma 9] that we get the right-hand side of

$$0 \leq \mathbb{E}_i[\varphi(q_i)] - \varphi(\mathbb{E}_i[q_i]) \leq D_\varphi \left( w \left\| (\varphi')^{-1} \left( \frac{\varphi(u) - \varphi(v)}{u - v} \right) \right\| \right), \quad (21)$$

478 where we can pick  $w \in \{u, v\}$  and  $D_\varphi$  is the Bregman divergence with generator  $\varphi$  (the left-hand  
 479 side is Jensen's inequality). Picking  $u \doteq -v$  (assuming wlog  $v > 0$ ) and letting

$$Y_v \doteq \frac{1-v}{2v} \cdot \log \left( 1 + \frac{2v}{1-v} \right), \quad (22)$$

480 for the choice  $\varphi(z) \doteq -\log(1+z)$  and  $w \doteq -v$ , we obtain after simplification

$$0 \leq \log \left( 1 + \frac{1}{n} \cdot \sum_i q_i \right) - \frac{1}{n} \cdot \sum_i \log(1+q_i) \leq Y_v - \log(Y_v) - 1. \quad (23)$$

481 Analyzing function  $v \mapsto Y_v - \log(Y_v) - 1$  reveals that it is upperbounded by  $v \mapsto v^2$  if  $v \in$   
 482  $[-8/10, 8/10]$ . Hence, multiplying by  $n$  all sides and passing to the exponential, we get that

$$1 \leq \frac{(1 + \frac{1}{n} \cdot \sum_i q_i)^n}{\prod_i (1 + q_i)} \leq \exp(nv^2), \forall 0 \leq v < \frac{8}{10},$$

483 which leads to the statement of the Lemma.  $\square$

484 Let us come back to our Riemannian summation setting and let

$$\begin{aligned} q_i &\doteq (1-t) \cdot |\mathbb{I}_i| f(\xi_i), \forall i \in [n], \\ v &\doteq \max_i |q_i|. \end{aligned}$$

485 Assume now that  $f$  is Riemann integrable, which allows to guarantee that, at least for  $n$  large enough  
 486 and a step of division  $\Delta_n$  not too large, we have  $|q_i| \leq 8/10, \forall i \in [n]$  and so we can use Lemma 2.

487 We get from (19) and Lemma 2, if  $t < 1$

$$\begin{aligned} S_{\Delta_n}^{(t)}(f) &= \frac{1}{1-t} \cdot \left( \prod_{i=1}^n (1+q_i) - 1 \right) \\ &\in \left[ \frac{1}{1-t} \cdot \left( \left( 1 + \frac{1}{n} \cdot \sum_i q_i \right)^n \cdot \exp(-nv \cdot v) - 1 \right), \frac{1}{1-t} \cdot \left( \left( 1 + \frac{1}{n} \cdot \sum_i q_i \right)^n - 1 \right) \right] \end{aligned} \quad (24)$$

488 and we permute the bounds if  $t > 1$ . Importantly, we note that

$$\sum_i q_i = (1-t) \cdot S_{\Delta_n}^{(1)}(f). \quad (25)$$

489 So suppose  $\lim_{n \rightarrow +\infty} S_{\Delta_n}^{(1)}(f) \doteq L_1 \doteq \int_a^b f(u) du$  is finite and choose the step of division  $\Delta_n$  not too  
490 large so that

$$n \cdot \max_i |q_i| = (1-t) \cdot \max_i (n|\mathbb{I}_i|) |f(\xi_i)| \leq K \cdot ((1-t)|b-a| \max f([a,b])),$$

491 for some constant  $K \geq 1$  (this is possible because  $f$  is (1-)Riemann integrable; we also note that if the  
492 division is regular then we can choose  $K = 1$ ). The value of  $K$  is not important: what is important is  
493 that  $nv$  remains finite in (24) as  $n$  increases, while  $\lim_{n \rightarrow +\infty} v = 0$ . Hence,  $\exp(-nv \cdot v) \rightarrow 1$  as  
494  $n \rightarrow +\infty$  and (24) implies, because  $\lim_{n \rightarrow +\infty} (1+a/n)^n = \exp a$ , the two first identities in

$$\begin{aligned} \lim_{n \rightarrow +\infty} S_{\Delta_n}^{(t)}(f) &= \frac{1}{1-t} \cdot \left( \lim_{n \rightarrow +\infty} \left( 1 + \frac{1}{n} \cdot \sum_i q_i \right)^n - 1 \right) \\ &= \frac{1}{1-t} \cdot \left( \exp \left( \sum_i q_i \right) - 1 \right) \\ &= \frac{1}{1-t} \cdot (\exp((1-t)L_1) - 1) \\ &= \log_t \exp(L_1) \\ &= \log_t \exp \left( \int_a^b f(u) du \right) \end{aligned} \quad (26)$$

495 ((26) follows from (25)) and by definition  $\lim_{n \rightarrow +\infty} S_{\Delta_n}^{(t)}(f) \doteq \int_a^b f(u) d_t u$ . So we get that  
496 Riemann integration ( $t = 1$ ) grants

$$\int_a^b f(u) d_t u = \log_t \exp \int_a^b f(u) du,$$

497 which, in addition to showing (14) (main file) also shows that  $t = 1$ -Riemann integration is equivalent  
498 to all  $t \neq 1$ -Riemann integration, ending the proof of Theorem 2.

499 **Remark 1.** *The absence of affine terms in upperbounding  $v \mapsto Y_v - \log(Y_v) - 1$  in the neighborhood  
500 of 0 is crucial to get to our result.*

501 **C.II Proof of Theorem 3**

502 Using Theorem 2, we just have to analyze the limit in relationship to the Riemannian case:

$$\begin{aligned} D_t F(z) &\doteq \lim_{\delta \rightarrow 0} \frac{\int_a^{z+\delta} f(u) d_t u \ominus_t \int_a^z f(u) d_t u}{\delta} \\ &= \lim_{\delta \rightarrow 0} \frac{\log_t \exp \int_a^{z+\delta} f(u) du \ominus_t \log_t \exp \int_a^z f(u) du}{\delta} \end{aligned} \quad (27)$$

$$= \lim_{\delta \rightarrow 0} \frac{1}{\delta} \cdot \log_t \left( \frac{\exp \int_a^{z+\delta} f(u) du}{\exp \int_a^z f(u) du} \right) \quad (28)$$

$$\begin{aligned} &= \lim_{\delta \rightarrow 0} \frac{1}{\delta} \cdot \log_t \exp \left( \int_a^{z+\delta} f(u) du - \int_a^z f(u) du \right) \\ &= \lim_{\delta \rightarrow 0} \frac{1}{\delta} \cdot \left( \int_a^{z+\delta} f(u) du - \int_a^z f(u) du \right) \end{aligned} \quad (29)$$

$$\doteq D_1 F(z) = f,$$

503 where the last identity is the classical Riemannian case, (27) follows from Theorem 2, (28) follows  
 504 from a property of  $\log_t (\log_t a \ominus \log_t b = \log_t(a/b))$ , (29) follows from the fact that  $\log_t \exp(z) =_0$   
 505  $z + o(z)$ . This completes the proof of Theorem 3.

506 **C.III Additional and helper results for Theorems 2 and 3**

507 We first state a trivial but important Lemma, whose content is partially used in the main file.

508 **Lemma 3.**  $g_t$  satisfies the following properties:

- 509 1.  $g_t(z)$  is strictly increasing for any  $t \in \mathbb{R}$ , strictly concave for any  $t > 1$ , strictly convex for  
 510 any  $t < 1$  and such that  $\text{sign}(g_t(z)) = \text{sign}(z), \forall t \in \mathbb{R}$ ;
- 511 2.  $D_t g_t(z) = 1$ ;
- 512 3. (t-integral mean-value) Suppose  $f$  Riemann integrable over  $[a, b]$ . Then there exists  $c \in$   
 513  $(a, b)$  such that

$$(b-a) \cdot g_{t'} \circ f(c) = \int_a^b f(u) d_t u \quad (t' \doteq 1 - (1-t)(b-a)).$$

514 We list a series of consequences to both theorems.

515 **General properties of t-integrals** Some properties generalize those for classical Riemann integration.

517 **Theorem C.1.** The following relationships hold for any  $t \in \mathbb{R}$  and any functions  $f, g$  t-Riemann  
 518 integrable over some interval  $[a, b]$ :

$$\begin{aligned} \int_a^b f(u) \{+ \text{ or } -\} g(u) d_t u &= \left( \int_a^b f(u) d_t u \right) \{ \oplus_t \text{ or } \ominus_t \} \left( \int_a^b g(u) d_t u \right) \quad (\text{additivity}), \\ \int_a^b \lambda f(u) d_t u &= \lambda \cdot \int_a^{(1-(1-t)\lambda)} f(u) d_t u \quad (\lambda \in \mathbb{R}) \quad (\text{dilativity}), \\ \int_a^b f(x) d_t u &= \int_a^c f(u) d_t u \oplus_t \int_c^b f(u) d_t u \quad (c \in [a, b]) \quad (\text{Chasles' relationship}), \\ \left| \int_a^b f(u) d_t u \right| &\leq \int_a^{(1-|1-t|)} |f(u)| d_t u \quad (\text{triangle inequality}), \\ \int_a^b f(u) d_t u &\leq \int_a^b g(u) d_t u \quad (f \leq g) \quad (\text{monotonicity}). \end{aligned}$$

519 *Proof.* We show additivity for  $\oplus_t/+$  (the same path shows the result for  $\ominus_t/-$ ):

$$\begin{aligned}
{}^{(t)} \int_a^b f(u) d_t u \oplus_t {}^{(t)} \int_a^b g(u) d_t u &= \log_t \exp \int_a^b f(u) du \oplus_t \log_t \exp \int_a^b g(u) du \\
&= \log_t \left( \exp \int_a^b f(u) du \cdot \exp \int_a^b g(u) du \right) \\
&= \log_t \exp \int_a^b (f+g)(u) du \\
&= {}^{(t)} \int_a^b (f+g)(u) d_t u.
\end{aligned}$$

520 We show dilativity, using  $t' \doteq 1 - (1-t)\lambda$ :

$$\begin{aligned}
{}^{(t)} \int_a^b \lambda f(u) d_t u &= \log_t \exp \int_a^b \lambda f(u) du \\
&= \log_t \exp \lambda \int_a^b f(u) du \\
&= \frac{1}{1-t} \cdot \left( \exp \left( \lambda(1-t) \int_a^b f(u) du \right) - 1 \right) \\
&= \frac{1-t'}{1-t} \cdot \frac{1}{1-t'} \cdot \left( \exp \left( (1-t') \int_a^b f(u) du \right) - 1 \right) \\
&= \lambda \cdot \frac{1}{1-t'} \cdot \left( \exp \left( (1-t') \int_a^b f(u) du \right) - 1 \right) \\
&= \lambda \cdot \log_{t'} \exp \int_a^b f(u) du \\
&= \lambda \cdot {}^{(t')} \int_a^b f(u) d_t u = \lambda \cdot {}^{(1-(1-t)\lambda)} \int_a^b f(u) d_t u.
\end{aligned}$$

521 The triangle inequality follows from the relationship:

$$|\log_t \exp(z)| \leq \log_{1-|1-t|} \exp |z|, \forall z, t \in \mathbb{R},$$

522 from which we use the fact that Riemann integration satisfies the triangle inequality and  $\log_{1-|1-t|}$  is  
523 monotonically increasing on  $\mathbb{R}_+$  in the penultimate line of:

$$\begin{aligned}
\left| {}^{(t)} \int_a^b f(u) d_t u \right| &\doteq \left| \log_t \exp \int_a^b f(u) du \right| \\
&\leq \log_{1-|1-t|} \exp \left| \int_a^b f(u) du \right| \\
&\leq \log_{1-|1-t|} \exp \int_a^b |f(u)| du \\
&= {}^{(1-|1-t|)} \int_a^b |f(u)| d_t u.
\end{aligned}$$

524 We show Chasles' relationship:

$$\begin{aligned}
{}^{(t)} \int_a^c f(u) d_t u \oplus_t {}^{(t)} \int_c^b f(u) d_t u &\doteq \log_t \exp \int_a^c f(u) du \oplus_t \log_t \exp \int_c^b f(u) du \\
&= \log_t \left( \exp \int_a^c f(u) du \cdot \exp \int_c^b f(u) du \right) \\
&= \log_t \exp \left( \int_a^c f(u) du + \int_c^b f(u) du \right) \\
&= \log_t \exp \int_a^b f(u) du \\
&= {}^{(t)} \int_a^b f(u) d_t u,
\end{aligned}$$

525 where the second identity uses the property  $\log_t a \oplus \log_t b = \log_t(ab)$ . Monotonicity follows  
526 immediately from the fact that  $z \mapsto \log_t \exp z$  is strictly increasing:

$$\begin{aligned}
{}^{(t)} \int_a^b f(u) d_t u &= \log_t \exp \int_a^b f(u) du \\
&\leq \log_t \exp \int_a^b g(u) du \doteq {}^{(t)} \int_a^b g(u) d_t u.
\end{aligned}$$

527 This ends the proof of Theorem C.1.  $\square$

528 **Computing  $t$ -integrals** Next, classical relationships to compute integrals do generalize to  $t$ -  
529 integration. We cite the case of integration by part.

530 **Lemma 4.** *Integration by part translates to  $t$ -integration by part as:*

$${}^{(t)} \int_a^b f g' du = {}^{(t)} [fg]_a^b \ominus_t {}^{(t)} \int_a^b f' g du,$$

531 where we let

$${}^{(t)} [h]_a^b \doteq \log_t \exp(h(b)) \ominus_t \log_t \exp(h(a)).$$

532 (Proof immediate from Theorem 2)

533 **Geometric properties based on  $t$ -integrals** This is a more specific result, important in the context  
534 of hyperbolic geometry: the well-known Hyperbolic Pythagorean Theorem (HPT) does translate  
535 to a tempered version with the same relationship to the Euclidean theorem. Consider a hyperbolic  
536 right triangle with hyperbolic lengths  $a, b, c, c$  being the hyperbolic length of the hypotenuse. Let  
537  $a_t, b_t, c_t$  denote the corresponding tempered lengths, which are therefore explicitly related using  $\mathfrak{g}_t$  as

$$a_t = \log_t \exp a, \quad b_t = \log_t \exp b, \quad c_t = \log_t \exp c.$$

538 Define the tempered generalization of cosh:

$$\cosh_t z \doteq \frac{\exp_t z + \exp_t(-z)}{2}. \tag{30}$$

539 The HPT tells us that  $\cosh c = \cosh a \cosh b$ . It is a simple matter of plugging  $\mathfrak{g}_t$ , using the fact that  
540  $\log_t$  and  $\exp_t$  are inverse of each other and simplifying to get the tempered HPT, which we call  
541  $t$ -HPT for short.

542 **Lemma 5.** ( *$t$ -HPT*) For any hyperbolic triangle described above, the tempered lengths are related as

$$\cosh_t c_t = \cosh_t a_t \cosh_t b_t.$$

543 Now, remark that for any  $t \neq 0$ , a series expansion around 0 gives

$$\exp_t(z) = 1 + \frac{tz^2}{2} + o(z^3)$$

544 ( $\exp_t$  is always infinitely differentiable around 0, for any  $t$ ) So the  $t$ -HPT gives

$$1 + \frac{tc_t^2}{2} + o(c_t^3) = \left(1 + \frac{ta_t^2}{2} + o(a_t^3)\right) \cdot \left(1 + \frac{tb_t^2}{2} + o(b_t^3)\right),$$

545 which simplifies, if we multiply both sides by  $2/t$  and simplify it into

$$c_t^2 + o(c_t^3) = a_t^2 + b_t^2 + o(a_t^3) + o(b_t^3),$$

546 which for an infinitesimal right triangle gives  $c_t^2 \approx a_t^2 + b_t^2$ , i.e. Pythagoras Theorem, as does the  
547 HPT one gives in this case ( $c^2 \approx a^2 + b^2$ ), which is equivalent of the particular  $t = 1$ -HPT case.

548  **$t$ -mean value Theorem** The  $t$ -derivative yields a generalization of the Euclidean mean-value  
549 theorem.

550 **Lemma 6.** Let  $t \in \mathbb{R}$  and  $f$  be continuous over an interval  $[a, b]$ , differentiable on  $(a, b)$  and such  
551 that  $-1/(1-t) \notin f([a, b])$ . Then  $\exists c \in (a, b)$  such that

$$D_t f(c) = \frac{(f(b) \ominus_t f(c)) - (f(a) \ominus_t f(c))}{b - a}.$$

552 *Proof.* We can obtain a direct expression of  $D_t f$  by using the definition of  $\ominus_t$  and the classical  
553 derivative:

$$D_t f(x) = \lim_{\delta \rightarrow 0} \frac{1}{\delta} \cdot \frac{f(x + \delta) - f(x)}{1 + (1-t)f(x)} = \frac{1}{1 + (1-t)f(x)} \cdot \lim_{\delta \rightarrow 0} \frac{f(x + \delta) - f(x)}{\delta} = \frac{f'(x)}{1 + (1-t)f(x)}. \quad (31)$$

554 From here, the mean-value theorem tells us that there exists  $c \in [a, b]$  such that

$$f'(c) = \frac{f(b) - f(a)}{b - a}$$

555 Dividing by  $1 + (1-t)f(c)$  (assuming  $f(c) \neq -1/(1-t)$ ) and reorganising, we get

$$\begin{aligned} D_t f(c) &= \frac{1}{b - a} \cdot \frac{f(b) - f(a)}{1 + (1-t)f(c)} \\ &= \frac{1}{b - a} \cdot \frac{f(b) - f(c)}{1 + (1-t)f(c)} - \frac{1}{b - a} \cdot \frac{f(a) - f(c)}{1 + (1-t)f(c)} \\ &= \frac{(f(b) \ominus_t f(c)) - (f(a) \ominus_t f(c))}{b - a}, \end{aligned}$$

556 which completes the proof of the Lemma.  $\square$

557 In fact, the  $t$ -derivative of  $f$  at some  $c$  is "just" an Euclidean derivative for an affine transformation of  
558 the function, namely  $z \mapsto f(z) \ominus_t f(c)$ , also taken at  $z = c$ . This "proximity" between  $t$ -derivation  
559 and derivation is found in the tempered chain rule (proof straightforward).

560 **Lemma 7.** Suppose  $g$  differentiable at  $z$  and  $f$  differentiable at  $g(z)$ . Then

$$D_t(f \circ g)(z) = D_t(f)(g(z)) \cdot g'(z).$$

#### 561 C.IV Sets endowed with a distortion and their $t$ -self: statistical information

562 Here,  $\mathcal{X}$  contains probability distributions or the parameters of probability distributions:  $f$  can then  
563 be an  $f$ -divergence (information theory) or a Bregman divergence (information geometry). Tsallis  
564 divergence and the tempered relative entropy are examples of  $t$ -additive information theoretic and  
565 information geometric divergences. A key property of information theory is the data processing  
566 inequality (**D**),  $\mathcal{X}$  being a probability space, which says that passing random variables through a  
567 Markov chain cannot increase their divergence as quantified by  $f$  [25, 35]. A key property of

568 information geometry is the population minimizer property **(P)**, which elicits a particular function of  
 569 a set of points as the minimizer of the expected distortion to the set, as quantified by  $f$  [3]. We let **(J)**  
 570 denote the joint convexity property, which would state for  $f$  and any  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2$  that

$$f(\lambda \cdot \mathbf{x}_1 + (1 - \lambda) \cdot \mathbf{x}_2, \lambda \cdot \mathbf{y}_1 + (1 - \lambda) \cdot \mathbf{y}_2) \leq \lambda f(\mathbf{x}_1, \mathbf{y}_1) + (1 - \lambda) f(\mathbf{x}_2, \mathbf{y}_2), \quad (32)$$

571 and convexity **(C)**, which amounts to picking  $\mathbf{y}_1 = \mathbf{y}_2$  (convexity in the left parameter) xor  $\mathbf{x}_1 = \mathbf{x}_2$   
 572 (in the right parameter).

573 **Lemma 8.** *For any  $t \in \mathbb{R}$ , the following holds true:*

574 **(D)**  *$f$  satisfies the data processing inequality iff  $f^{(t)}$  satisfies the data processing inequality;*

575 **(P)**  $\mu_* \in \arg \min_{\mu} \sum_i f(\mathbf{x}_i, \mu)$  iff

$$\mu_* \in \arg \min_{\mu} (\oplus_t)_i f^{(t)}(\mathbf{x}_i, \mu); \quad (33)$$

576 **(J)**  *$f$  satisfies (32) iff the following  $(t, t', t'')$ -joint convexity property holds:*

$$f^{(t)}(\lambda \cdot \mathbf{x}_1 + (1 - \lambda) \cdot \mathbf{x}_2, \lambda \cdot \mathbf{y}_1 + (1 - \lambda) \cdot \mathbf{y}_2) \leq \lambda f^{(t')}(x_1, y_1) + (1 - \lambda) f^{(t'')}(x_2, y_2), \quad (34)$$

577 *with  $t' \doteq \min\{t, 1 - \lambda + \lambda t\}$  and  $t'' \doteq \min\{t, \lambda + (1 - \lambda)t\}$ .*

578 *Proof.* **(D)** and **(P)** are immediate consequences of Lemma 3 (point [1.], main file) and properties of  
 579  $\log_t$ . We prove **(J)**.  $g_t$  being strictly increasing for any  $t$ , we get for  $t \leq 1$

$$\begin{aligned} f^{(t)}(\lambda \cdot \mathbf{x}_1 + (1 - \lambda) \cdot \mathbf{x}_2, \lambda \cdot \mathbf{y}_1 + (1 - \lambda) \cdot \mathbf{y}_2) &\leq g_t(\lambda f(\mathbf{x}_1, \mathbf{y}_1) + (1 - \lambda) f(\mathbf{x}_2, \mathbf{y}_2)) \\ &\leq \lambda \cdot g_t \circ f(\mathbf{x}_1, \mathbf{y}_1) + (1 - \lambda) \cdot g_t \circ f(\mathbf{x}_2, \mathbf{y}_2) \\ &= \lambda f^{(t)}(\mathbf{x}_1, \mathbf{y}_1) + (1 - \lambda) f^{(t)}(\mathbf{x}_2, \mathbf{y}_2), \end{aligned} \quad (35)$$

580 because  $g_t$  is convex. If  $t > 1$ , we restart from the first inequality and remark that

$$\begin{aligned} g_t(\lambda f(\mathbf{x}_1, \mathbf{y}_1) + (1 - \lambda) f(\mathbf{x}_2, \mathbf{y}_2)) &= g_t(\lambda f(\mathbf{x}_1, \mathbf{y}_1)) \oplus_t g_t((1 - \lambda) f(\mathbf{x}_2, \mathbf{y}_2)) \\ &\leq g_t(\lambda f(\mathbf{x}_1, \mathbf{y}_1)) + g_t((1 - \lambda) f(\mathbf{x}_2, \mathbf{y}_2)) \\ &= \lambda \cdot g_{1-\lambda+\lambda t} \circ f(\mathbf{x}_1, \mathbf{y}_1) + (1 - \lambda) \cdot g_{\lambda+(1-\lambda)t} \circ f(\mathbf{x}_2, \mathbf{y}_2) \\ &= \lambda f^{(1-\lambda+\lambda t)}(\mathbf{x}_1, \mathbf{y}_1) + (1 - \lambda) f^{(\lambda+(1-\lambda)t)}(\mathbf{x}_2, \mathbf{y}_2). \end{aligned} \quad (36)$$

581 The inequality is due to the fact that  $a \oplus_t b = a + b + (1 - t)ab \leq a + b$  if  $ab \geq 0$  and  $t \geq 1$ . The  
 582 last equality holds because, for  $t' \doteq 1 - (1 - t)b$ , we have

$$\log_t(a^b) = \frac{a^{(1-t)b} - 1}{1 - t} = \frac{1 - t'}{1 - t} \cdot \frac{a^{1-t'} - 1}{1 - t'} = b \cdot \log_{t'}(a).$$

583 Putting altogether (35) and (36), we get that for any  $t \in \mathbb{R}$ ,

$$\begin{aligned} f^{(t)}(\lambda \cdot \mathbf{x}_1 + (1 - \lambda) \cdot \mathbf{x}_2, \lambda \cdot \mathbf{y}_1 + (1 - \lambda) \cdot \mathbf{y}_2) &\leq \lambda f^{(\min\{t, 1 - \lambda + \lambda t\})}(\mathbf{x}_1, \mathbf{y}_1) + (1 - \lambda) f^{(\min\{t, \lambda + (1 - \lambda)t\})}(\mathbf{x}_2, \mathbf{y}_2), \end{aligned} \quad (37)$$

584 as claimed for **(J)**. This ends the proof of Lemma 8.  $\square$

585 We note that (34) also translates into a property for **(C)**; also, if  $t \leq 1$ , then  $t = t' = t''$  in (34).

## 586 C.V Using the T-calculus in the Lorentz model of hyperbolic geometry

587 As an additional example of use, we consider the Lorentz model for a simple illustration in which  
 588 one creates approximate metricity in the t-self. This  $d$ -dimensional manifold is embedded in  $\mathbb{R}^{d+1}$   
 589 via the hyperboloid with constant  $-c < 0$  curvature and defined by  $\mathbb{H}_c \doteq \{\mathbf{x} \in \mathbb{R}^{d+1} : x_0 >$   
 590  $0 \wedge \mathbf{x} \circ \mathbf{x} = -1/c\}$ , with  $\mathbf{x} \circ \mathbf{y} \doteq -x_0 y_0 + \sum_{i=1}^d x_i y_i$  the Lorentzian inner product [27, Chapter  
 591 3]. In ML, two distortions are considered on  $\mathbb{H}_c$  [17], one of which is the Lorentzian "distance",  
 592  $d_L(\mathbf{x}, \mathbf{y}) \doteq -(2/c) - 2 \cdot \mathbf{x} \circ \mathbf{y}$ . It is notoriously not a distance because it does not satisfy **(T)**, yet we  
 593 show that we can pick  $t$  and the curvature in such a way that the t-self is arbitrarily close to a metric  
 594 space.

595 **Lemma 9.**  $\forall \delta > 0$ , pick  $t$  and curvature  $c$  as  $t = 1 + (1/\delta)$ ,  $c = 2/\delta$ . Then the  $t$ -self of  $\mathbb{H}_c$  is  
 596 approximately metric:  $d_L^{(t)}$  satisfies **(R)**, **(S)**, **(I)**, and the  $\delta$ -triangle inequality,

$$d_L^{(t)}(\mathbf{x}, \mathbf{z}) \leq d_L^{(t)}(\mathbf{x}, \mathbf{y}) + d_L^{(t)}(\mathbf{y}, \mathbf{z}) + \delta, \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{H}_c. \quad (38)$$

597 *Proof.*  $d_L^{(t)}$  still obviously satisfies reflexivity, the identity of indiscernibles, and non-negativity, so  
 598 we check the additional property it now satisfies, the weaker version of the triangle inequality. Given  
 599 any  $\mathbf{x}, \mathbf{y}, \mathbf{z}$  in  $\mathbb{H}_c$ , condition  $d_L^{(t)}(\mathbf{x}, \mathbf{z}) \leq d_L^{(t)}(\mathbf{x}, \mathbf{y}) + d_L^{(t)}(\mathbf{y}, \mathbf{z}) + \delta$  for  $t > 1$  is

$$\frac{1 - \exp(-(t-1)(-\frac{2}{c} - 2 \cdot \mathbf{x} \circ \mathbf{z}))}{t-1} \leq \frac{1 - \exp(-(t-1)(-\frac{2}{c} - 2 \cdot \mathbf{x} \circ \mathbf{y}))}{t-1} + \frac{1 - \exp(-(t-1)(-\frac{2}{c} - 2 \cdot \mathbf{y} \circ \mathbf{z}))}{t-1} + \delta,$$

600 which simplifies to

$$\begin{aligned} \exp(2(t-1) \cdot \mathbf{x} \circ \mathbf{y}) + \exp(2(t-1) \cdot \mathbf{y} \circ \mathbf{z}) &\leq \exp\left(-\frac{2(t-1)}{c}\right) + (t-1)\delta \exp\left(-\frac{2(t-1)}{c}\right) \\ &\quad + \exp(2(t-1) \cdot \mathbf{x} \circ \mathbf{z}). \end{aligned}$$

601 We have  $\mathbf{x} \circ \mathbf{y} \leq -1/c$  by definition, so a sufficient condition to get the inequality is to have

$$\exp(2(t-1) \cdot \mathbf{y} \circ \mathbf{z}) \leq (t-1)\delta \exp\left(-\frac{2(t-1)}{c}\right). \quad (39)$$

602 Function  $h(z) \doteq z\delta \exp(-2z/c)$  is maximum over  $\mathbb{R}_+$  for  $z_* = c/2$ , for which it equals  $h(z_*) =$   
 603  $c\delta/(2e)$ . Fix  $t = 1 + (c/2)$ . We then have  $\exp(2(t-1) \cdot \mathbf{y} \circ \mathbf{z}) \leq 1/e$  so to get (39) for this choice  
 604 of  $t$ , it is sufficient to pick curvature  $c = 2/\delta$ , yielding relationship  $t = 1 + (1/\delta)$ .  $\square$

## 605 C.VI Proof of Lemma 1

606 We start by proving condition

$$d_*^{(t)} \geq g(k) \quad (40)$$

607 in the Lemma. Using the proof of **(author?)** [19, Proposition 3.1], we know that any point  $\mathbf{x}$   $k$ -close  
 608 to the boundary (Definition 5.1) satisfies

$$d^{(t)}(\mathbf{x}, \mathbf{0}) = \log_t \left( \frac{1 + \|\mathbf{x}\|}{1 - \|\mathbf{x}\|} \right) = \log_t (2 \cdot 10^k - 1),$$

609 so to get (40), we want  $\log_t(2 \cdot 10^k - 1) \geq g(k)$ . Letting  $t = 1 - f(k)$  with  $f(k) \in \mathbb{R}$ , we observe  
 610  $\log_t(2 \cdot 10^k - 1) = ((2 \cdot 10^k - 1)^{f(k)} - 1)/f(k)$ , so we want, after taking logs,

$$\log(2 \cdot 10^k - 1) \geq \frac{\log(1 + f(k)g(k))}{f(k)} \quad (41)$$

611 (this also holds if  $f(k) < 0$  because  $1 + f(k)g(k) < 1$ ), and there remains to observe  $\log(2 \cdot 10^k - 1) =$   
 612  $k \log(10) + \log(2 - 1/10^k)$  with  $\log(2 - 1/10^k) \geq 0, \forall k \geq 0$ . Hence, to get (41) it is sufficient to  
 613 request

$$\frac{\log(1 + f(k)g(k))}{f(k)} \leq \log(10) \cdot k,$$

614 which is (40). For such  $t$ , the new hyperbolic constant satisfies

$$\tau_t = \log_t \exp \tau = \frac{1}{f(k)} \cdot (\exp(f(k)\tau) - 1)$$

## 615 C.VII Proof of Theorem 1

616 Take any observation  $\mathbf{x} \in \mathcal{X}$ . Trivially,  $H'(\mathbf{x})$  is a real value that belongs to the path followed by  
 617  $\mathbf{x}$  in  $H$ . Furthermore, if we consider any path from the root to a leaf in  $H$ , all the vertices in its  
 618 strictly monotonically increasing subsequence of absolute confidences appear in  $H'$  (conditions 4, 8  
 619 in GETMDT), which guarantees the condition on prediction in **(M)**. Hence **(M)** is satisfied.

**Algorithm 2** LOGISTICBOOST( $\mathcal{S}$ )

---

**Input:** Labeled sample  $\mathcal{S} \doteq \{(\mathbf{x}_i, y_i), i \in [m]\}$ ,  $T \in \mathbb{N}_{>0}$ ;  
 1 :  $w_{1i} \leftarrow 1/2, \forall i \in [m]$ ; // initialize all weights (equivalent to maximally unconfident prediction)  
 2 :   **for**  $j = 1, 2, \dots, T$   
 3 :      $H_j \leftarrow \text{WEAK-LEARN}(\mathcal{S}, \mathbf{w}_j)$ ; // request a DT as a "weak hypothesis"  
 4 :      $\alpha_j \in \mathbb{R}$ ; // picks leveraging coefficient for  $H_j$   
 5 :     **for**  $i = 1, 2, \dots, m$  // weight update, not normalized

$$w_{(j+1)i} \leftarrow \frac{w_{ji}}{w_{ji} + (1 - w_{ji}) \cdot \exp(\alpha_j y_i H_j(\mathbf{x}_i))}; \quad (42)$$

**Output:** Classifier  $\mathbf{H}_T \doteq \sum_j \alpha_j H_j(\cdot)$ ;

---

621 We want a boosting algorithm for the liner combination of DTs which displays classification that can  
 622 be easily and directly embedded in the Poincaré disk. Algorithm LOGISTICBOOST is provided above.  
 623 For the weight update, we refer to [23]. We already know how to embed DTs via their Monotonic  
 624 DTs. What a boosting algorithm of this kind does is craft

$$\mathbf{H}_T \doteq \sum_{j=1}^T u_j(\cdot), \quad u_j(\mathbf{x}) \doteq \alpha_j \cdot H_j(\mathbf{x}); \quad (43)$$

625 Remark that we have merged the leveraging coefficient and the DTs' outputs  $H_j$ , on purpose. To  
 626 compute the leveraging coefficients  $\alpha_j$  in Step 4, we use AdaBoost's secant approximation trick<sup>†</sup>,  
 627 applied not to the exponential loss but to the logistic loss: for any  $z \in [-R, R]$  and  $\alpha \in \mathbb{R}$ ,

$$\log(1 + \exp(-\alpha z)) \leq \frac{1+u}{2} \cdot \log(1 + \exp(-\alpha R)) + \frac{1-u}{2} \cdot \log(1 + \exp(\alpha R)). \quad (44)$$

628 Hence, to compute the leveraging coefficient  $\alpha_j^*$  that approximately minimizes the current loss,  
 629  $\sum_i w_{ji} \log(1 + \exp(-\alpha_j y_i H_j(\mathbf{x}_i)))$ , we minimize instead the upperbound using (44). Letting

$$(\psi_{\text{LOG}})_j^* \doteq \max_{\lambda \in \Lambda(H_j)} \left| \log \left( \frac{p_{j\lambda}}{1 - p_{j\lambda}} \right) \right| \quad (45)$$

630 (we remind that  $\Lambda(\cdot)$  denotes the set of leaves of a tree; the index in the local proportion of positive  
 631 example  $p_{j\lambda}^+$  reminds that weights used need to be boosting's weights) which we note can be directly  
 632 approximated on the Poincaré disk by looking at the leaf nearest to the border of the disk. Because  
 633 the maximal absolute confidence in the DT is also the maximal absolute confidence in its MDT, we  
 634 obtain the sought minimum,

$$\alpha_j^* = \frac{1}{(\psi_{\text{LOG}})_j^*} \cdot \log \left( \frac{1+r_j}{1-r_j} \right),$$

635 where  $r_j \in [-1, 1]$  is the normalized edge

$$\begin{aligned} r_j &\doteq \frac{1}{\sum_i w_{ji}} \cdot \sum_{i \in [m]} w_{ji} \cdot \frac{y_i H_j(\mathbf{x}_i)}{\max_k |H_j(\mathbf{x}_k)|} \\ &= \mathbb{E}_{i \sim \tilde{w}_j} \left[ \frac{1}{(\psi_{\text{LOG}})_j^*} \cdot \log \left( \frac{p_{j\lambda(\mathbf{x}_i)}^+}{1 - p_{j\lambda(\mathbf{x}_i)}^+} \right) \right], \end{aligned} \quad (46)$$

636 where  $\tilde{w}_j$  indicates normalized weights. It is not hard to show that because we use the local posterior  
 637  $p_{j\lambda}^+$  at each leaf,  $\alpha_j^* \geq 0$  and also  $r_j \geq 0$ . Hence, everything is like if we had an imaginary node  $\nu_j$   
 638 with  $p_{j\nu_j}^+ \doteq (1+r_j)/2 (\geq 1/2)$  and positive confidence

$$(\psi_{\text{LOG}})_j \doteq \psi_{\text{LOG}}(p_{j\nu_j}^+) = \log \left( \frac{p_{j\nu_j}^+}{1 - p_{j\nu_j}^+} \right)$$

<sup>†</sup>Explained in (author?) [31, Section 3.1].

639 that we can display in Poincaré disk (we choose to do it as a circle; see Figure 2, main file). We  
 640 deduce from (43) that

$$u_j(\mathbf{x}) = \frac{(\psi_{\text{LOG}})_j}{(\psi_{\text{LOG}})_j^*} \cdot \log \left( \frac{p_{j\lambda(\mathbf{x})}^+}{1 - p_{j\lambda(\mathbf{x})}^+} \right),$$

641 and note that *all* three key parameters can easily be displayed or computed directly from the Poincaré  
 642 disk.

### 643 C.IX Modifying Sarkar’s embedding in Poincaré disk

644 Sarkar’s algorithm [30] gives a clean low-distortion embedding when the tree is binary or the arc  
 645 length is constant [29]. Things are different in our case: MDT nodes have arbitrary out-degrees,  
 646 and lengths depend on the absolute confidence at the corresponding MDT nodes. Plus, a direct  
 647 implementation of Sarkar’s algorithm would violate the constraint for strict path-monotonicity in an  
 648 MDT that nodes in a path from the root need to progressively come closer to the border – equivalently,  
 649 it would create substantial embedding errors for confidences and violate **(B)**. Without focusing on  
 650 an optimal solution (that we leave for future work), one can remark that all these problems can be  
 651 heuristically addressed by changing one step of Sarkar’s algorithm, replacing the use of the total  $2\pi$   
 652 fan out angle for mapping children (Step 5: in Algorithm 1 of [29]) by a variable angle with a special  
 653 orientation in the disk.

654 We refer to the concise and neat description of Sarkar’s embedding in [29] for the full algorithm. Our  
 655 modification relies on changing one step of the algorithm, as described in Figure I. The key step  
 656 that we change is step 5: in the description of **(author?)** [29, Algorithm 1]. This step embeds the  
 657 children of a given node (and then the algorithm proceeds recursively until all nodes are processed).  
 658 Sarkar’s algorithm corresponds to the simple case where all arcs to/from a node define a fixed angle,  
 659 which does not change during reflection because Poincaré model is conformal. Hence, if the tree is  
 660 binary, this angle is  $2\pi/3$ , which provides a very clean display of the tree. In our case however, some  
 661 children may have just one arc to a leaf while others may support big subtrees. Also, arc lengths  
 662 can vary substantially. We thus design the region of the disk into which the subtrees are going to  
 663 be embedded by choosing an angle proportional to the number of leaves reachable from the node,  
 664 and of course lengths have to match the difference between absolute confidence between a node  
 665 and its children. There is no optimization step to learn a clean embedding, so we rely on a set of  
 666 hyperparameters to effectively compute this new step of the algorithm.

## 667 D Supplementary material on experiments

### 668 D.I Domains

669 The domains we consider are all public domains, from the UCI repository of ML datasets [12],  
 670 OpenML, or Kaggle, see Table I.

### 671 D.II Visualizing a DT via its MDT

672 All test errors are estimated from a 10-fold stratified cross-validation. One can always choose to  
 673 directly learn Monotonic DTs instead of DTs, given their natural fit for embedding in the Poincaré  
 674 disk. In this case, the hyperbolic representation of the MDT can be immediately used for assessment.  
 675 Suppose we stick to learning DTs (because *e.g.* they have been standard in ML for decades) *and* wish  
 676 to use the visualization of its corresponding MDT to make inferences on the original DT itself. Can  
 677 we make reliable conclusions? We explore this question by observing that the prediction from DT to  
 678 MDT can only change if the leaf  $\lambda$  that an example reaches in the DT satisfies two conditions: (a)  
 679  $\lambda$  does not appear as a leaf in the MDT, and (b) it is tagged to an internal node of the MDT with  
 680 a confidence of the opposite sign (this is **(D1)**, Step 3: in GETMDT). Without tackling the formal  
 681 aspect of this question, we have designed a simple experiment for a simple assessment of whether  
 682 / when this is reasonable. We have trained a set of  $T=200$  boosted DTs, each with maximal size  
 683 200 (total number of nodes). After having computed the corresponding MDTs, we have compared  
 684 the test errors of the boosted set of trees and that of the ensemble where each DT is replaced by its  
 685 MDT, *but* the leveraging coefficients do not change. Intuitively, if test errors are on par, the variation

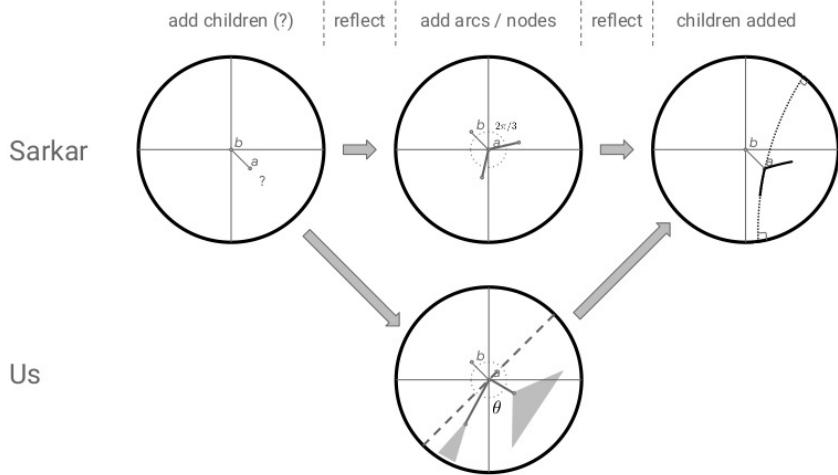


Figure I: Schematic description of our modification (bottom) of Sarkar’s algorithm (up). Our modification solely changes the step in which the children of a node (here,  $a$ ) are computed, this node having been reflected back to the origin. Instead of using a fixed angle and length to position arcs (and thus children of the node), the angle depends on the number of leaves that can be reached from a child, and the length depends on the difference between absolute confidence between the node and the corresponding child. We also define an angle which represents the domain (before reflecting back) in which the embedding is going to take place, shown with the thick dashed line (here, this angle is  $\pi$ ).

Domain	$m$	$d$	License
breastwisc	699	9	CC BY 4.0
ionosphere	351	33	CC BY 4.0
tictactoe	958	9	PDDL
winered	1 599	11	CC BY 4.0
german	1 000	20	CC BY 4.0
analcatdata_supreme	4 053	8	CC BY 4.0
abalone	4 177	8	CC BY 4.0
qsar	1 055	41	CC BY 4.0
hillnoise	1 212	100	CC BY 4.0
firmteacher	10 800	16	CC BY 4.0
online_shoppers_intention	12 330	17	CC BY 4.0
give_me_some_credit	120 269	11	None
Buzz_in_social_media (Tom’s hardware)	28 179	96	CC BY 4.0
Buzz_in_social_media (twitter)	583 250	78	CC BY 4.0

Table I: UCI, OpenML (Analcatdata\_supreme) and Kaggle (Give\_me\_some\_credit) domains considered in our experiments ( $m$  = total number of examples,  $d$  = number of features), ordered in increasing  $m \times n$ . Dataset licenses listed in last column.

in classification of DTs vs MDTs (including confidences) is negligible, and we can “reduce” the interpretation of the DT to that of its MDT in the Poincaré disk. The results are summarized in Table II.

From Table II, we can safely say that our hypothesis reasonably holds in many cases, with two important domains for which it does not: `ionosphere` and `hillnoise`. For the former domain, we attribute it to the small size of the domain, which prevents training big enough trees; for the latter domain, we attribute it to the fact that the domain contains substantial noise, which makes it difficult to substantially improve posteriors by splitting and thus make many DT nodes, including leaves, “disappear” in the MDT conversion (Steps 2: and 14: in GETMDT).

domain	DT	MDT	p-val
breastwisc	$4.15 \pm 2.18$	$5.00 \pm 2.38$	<b>0.2408</b>
ionosphere	$5.71 \pm 2.70$	$9.11 \pm 5.34$	0.0237
tictactoe	$2.61 \pm 1.85$	$2.09 \pm 1.10$	<b>0.1390</b>
winered	$18.39 \pm 2.02$	$18.32 \pm 2.27$	<b>0.9227</b>
german	$24.00 \pm 4.37$	$23.90 \pm 4.25$	<b>0.8793</b>
analcatdata_supreme	$23.40 \pm 1.85$	$22.56 \pm 1.75$	0.0134
abalone	$22.15 \pm 2.20$	$21.14 \pm 2.44$	<b>0.1267</b>
qsar	$12.98 \pm 2.63$	$12.89 \pm 4.05$	<b>0.8772</b>
hillnoise	$37.04 \pm 2.73$	$45.88 \pm 4.72$	0.0008
firmteacher	$6.87 \pm 1.23$	$7.04 \pm 1.17$	<b>0.4292</b>
online_shoppers_intention	$9.90 \pm 0.78$	$10.67 \pm 0.54$	0.0057
hardware	$1.33 \pm 0.27$	$1.44 \pm 0.19$	<b>0.0534</b>

Table II: Results of the experiments checking whether "reducing" the interpretation of a DT to that of its MDT (using its hyperbolic embedding) can give accurate information about the tree as well. Numerical columns, from left to right, give the average  $\pm$  std dev error for DTs and MDTs and provide the p-value of a Student paired t-test with H0 being the identity of the average errors. Entries in **bold faces** correspond to *keeping H0* for a 0.05 first-order risk. See text for details.

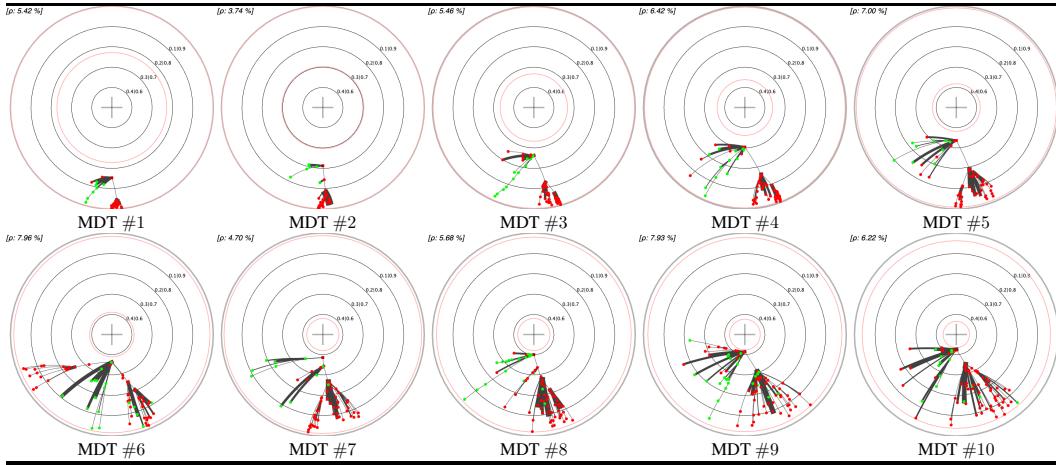


Table III: First 10 Monotonic Decision Trees (MDTs) embedded in Poincaré disk, corresponding to several Decision Trees (DTs) with  $\leq 200$  nodes each, learned by boosting the log / logistic-loss on UCI `online_shoppers_intention`. Isolines correspond to the node's prediction confidences, and geometric embedding errors ( $\rho\%$ ) are indicated. See text for details.

### 695 D.III All Poincaré disk embeddings

696 They are presented in Table III through to XI, showing the first models learned in a fold of the  
697 cross-validation experiments. producing Table II

### 698 D.IV Interest of the t-self for visualization

699 We now test the experimental impact of switching to the t-self in Poincaré disk model (See Table XII  
700 for a clear depiction of how changing  $t$  can yields better visualization close to the border of the disk).  
701 Recall that switching to the t-self moves way model parts close to  $\partial\mathbb{B}_1$  but keeps a low non-linear  
702 distortion around the center, which is thus roughly only affected by a scaling factor. Figure I presents  
703 a few results. For domain `online_shoppers_intention`, we note that the part of the tree that is  
704 within the isoline defined by posterior  $p^+ \in \{0.1, 0.9\}$  gets indeed just scaled: both plots look quite  
705 identical. Very substantial differences appear near the border: the best parts of the model could easily  
706 be misjudged as equivalent from  $\mathbb{B}_1$  alone (orange rectangle) but there is little double from  $\mathbb{B}_1^{(t)}$  that  
707 one of them, which crosses the  $p^+ \in \{0.001, 0.009\}$  isoline, is in fact much better than the others.

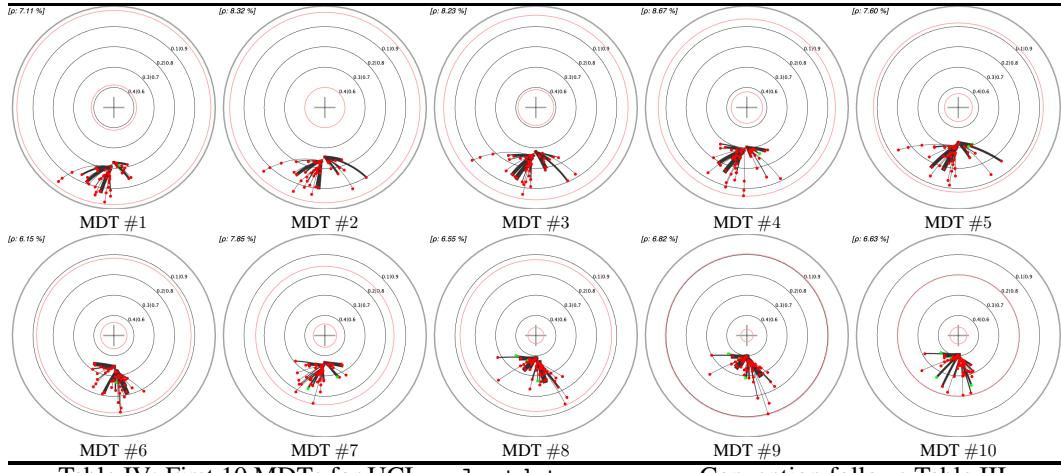


Table IV: First 10 MDTs for UCI analcatdata\_supreme. Convention follows Table III.

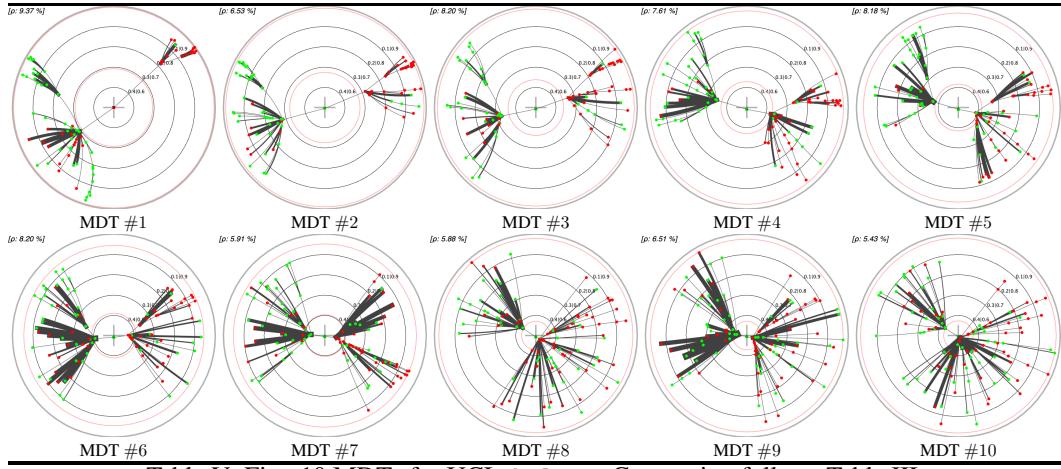


Table V: First 10 MDTs for UCI abalone. Convention follows Table III.

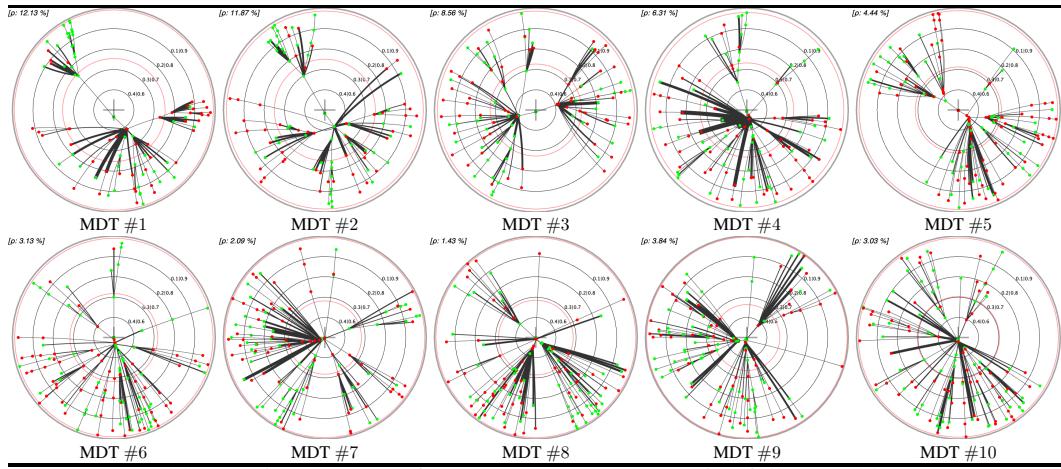


Table VI: First 10 MDTs for UCI winered. Convention follows Table III.

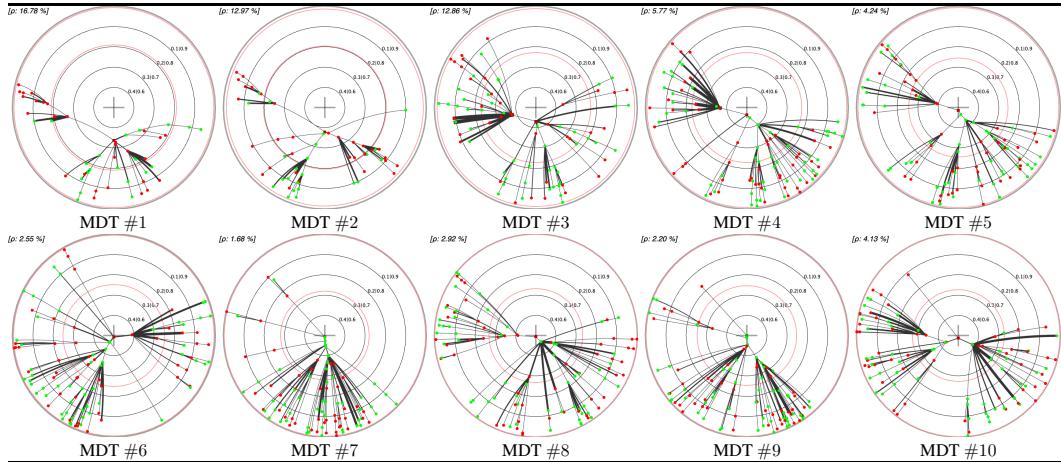


Table VII: First 10 MDTs for UCI qsar. Convention follows Table III.

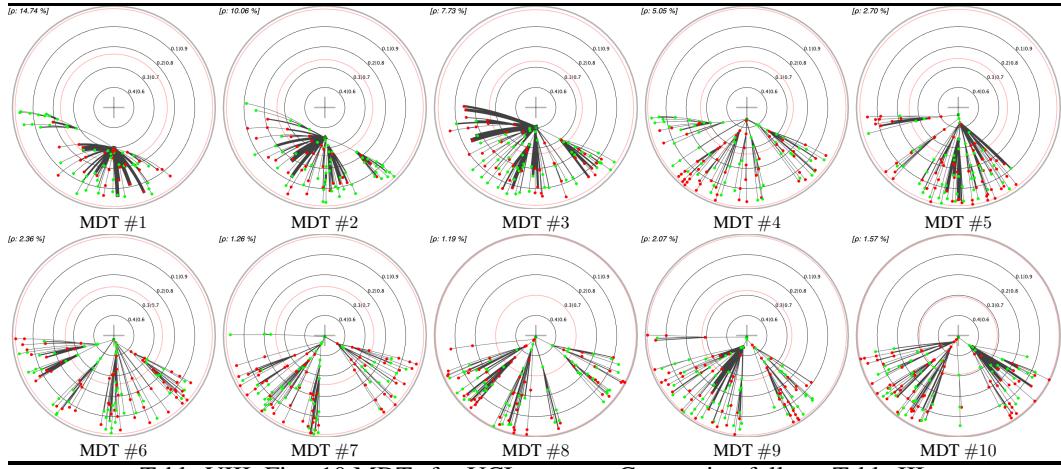


Table VIII: First 10 MDTs for UCI german. Convention follows Table III.

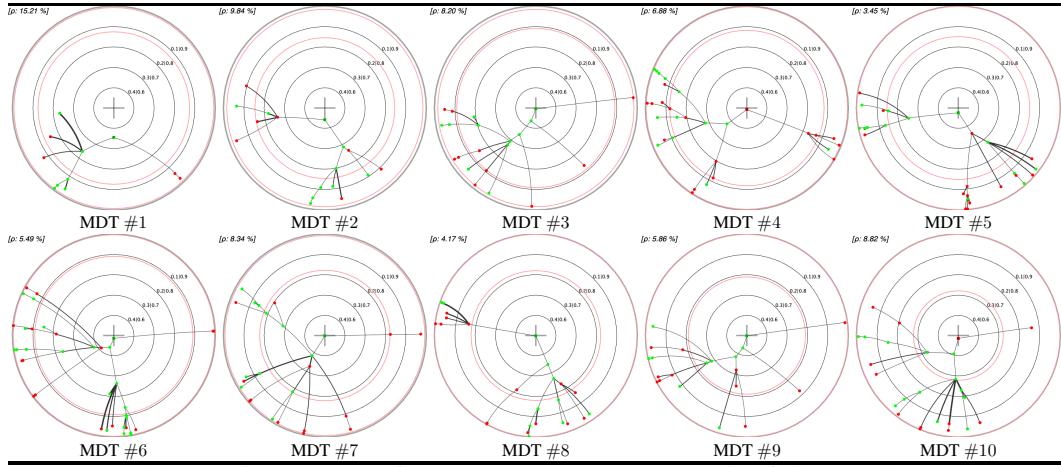


Table IX: First 10 MDTs for UCI ionosphere. Convention follows Table III.

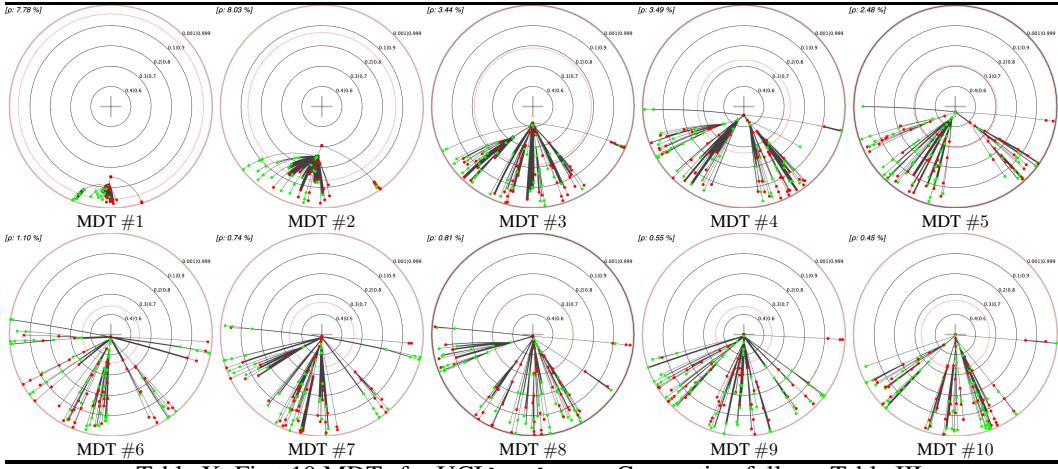


Table X: First 10 MDTs for UCI hardware. Convention follows Table III.

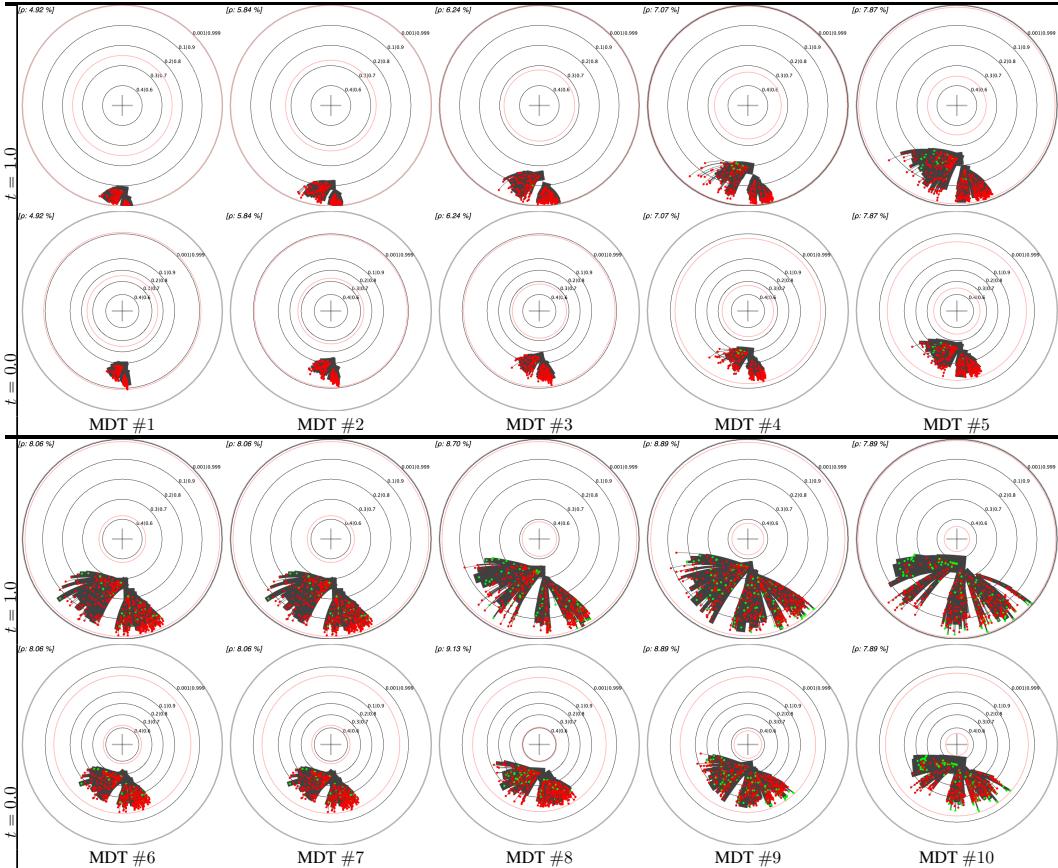


Table XI: Plot of the 10 first MDTs learned on kaggle `give_me_some_credit` (top panel and bottom panel). In each panel, we plot the embedding in  $\mathbb{B}_1$  (top row) and the t-self  $\mathbb{B}_1^{(0)}$  ( $t = 0$ , bottom row). Remark the ability for the t-self to display a clear difference between the best subtrees, subtrees that otherwise appear quite equivalent in terms of confidence from  $\mathbb{B}_1$  alone.

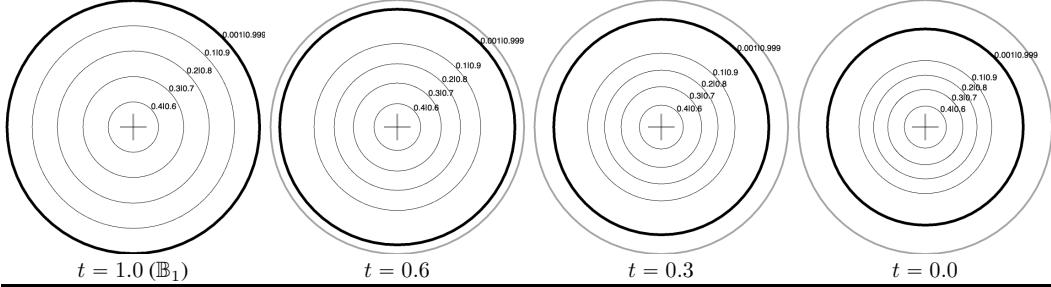


Table XII: T-self  $\mathbb{B}_1^{(t)}$  for Poincaré disk  $\mathbb{B}_1$  (left), for several  $t$ s. We plot the same set of isolines of  $\mathbb{B}_1$  (and their mapping via  $\mathfrak{p}_t$ ), parameterized by probability  $p \in [0, 1]$  which gives the norm  $r \doteq |2p - 1|$  in  $\mathbb{B}_1$  (From the innermost to the outermost,  $p \in \{0.6, 0.7, 0.8, 0.9, 0.999\}$ , or by symmetry for  $p' \doteq 1 - p$ , also indicated). Remark the equidistance of isolines in  $\mathbb{B}_1$ , approximately kept in  $\mathbb{B}_1^{(t)}$  for  $p$  up to 0.9 ( $p' = 0.1$ ), while the distortion we need clearly happens near the border: outermost isoline  $p \in \{0.001, 0.999\}$  (plotted with **bigger** width) is smoothly and substantially “moved” within the t-self as  $t$  decreases, guaranteeing good readability and coding convenience (see text).

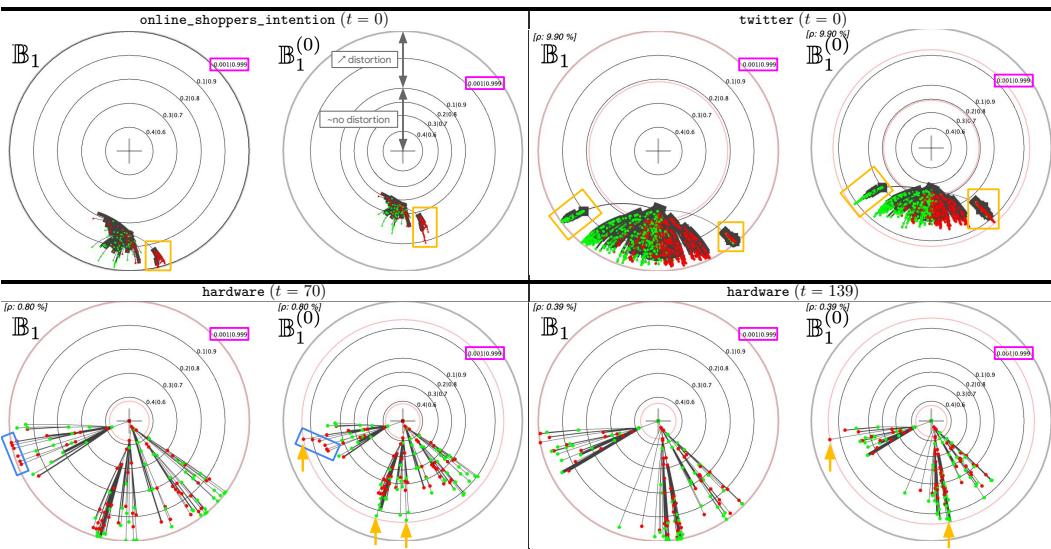


Figure I: *Top pane*: comparison between Poincaré disk embedding and its t-self for  $t = 0$  for an MDT learned on UCI `online_shoppers_intention` (left) and `twitter` (right). On the left panel, we have not plotted the boosting coefficients information. The isoline distinguished in **magenta** is the big **width** one in Table XII. Note the difference in (non-linear) distortion created in the t-self, in which the central part just enjoys a scaling. *Bottom pane*: stark differences in visualization between  $\mathbb{B}_1$  and the t-self do not just appear initially: they can appear at any iteration. We display two MDTs learned on UCI `hardware` at two different iterations (indicated). It would be very hard to differentiate the best leaves from just the Poincaré disk embedding, while it becomes obvious from the t-self (**orange** arrows). Note also the set of **red** nodes in the Poincaré disk for  $t = 70$  that mistakenly look aligned, but not in the t-self (**blue** rectangle). See text for details.

708 When many subtrees seem to be aggregating near the border as in `buzz_in_social_media`, stark  
 709 differences can appear on the t-self: the best subtrees are immediately spotted from the t-self (**orange**  
 710 rectangles). In between, the t-self makes a much more visible ordering between the best nodes and  
 711 subtrees, compared to Poincaré disk. `hardware` demonstrate that such very good nodes that are hard  
 712 to differentiate from the others in  $\mathbb{B}_1$  can appear at any iteration.

713 **NeurIPS Paper Checklist**

714 The checklist is designed to encourage best practices for responsible machine learning research,  
715 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
716 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
717 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
718 towards the page limit.

719 Please read the checklist guidelines carefully for information on how to answer these questions. For  
720 each question in the checklist:

- 721 • You should answer [Yes] , [No] , or [NA] .
- 722 • [NA] means either that the question is Not Applicable for that particular paper or the  
723 relevant information is Not Available.
- 724 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

725 **The checklist answers are an integral part of your paper submission.** They are visible to the  
726 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it  
727 (after eventual revisions) with the final version of your paper, and its final version will be published  
728 with the paper.

729 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
730 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a  
731 proper justification is given (e.g., "error bars are not reported because it would be too computationally  
732 expensive" or "we were unable to find the license for the dataset we used"). In general, answering  
733 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we  
734 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
735 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
736 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
737 please point to the section(s) where related material for the question can be found.

738 **IMPORTANT**, please:

- 739 • **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”,**
- 740 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 741 • **Do not modify the questions and only use the provided macros for your answers.**

742 **1. Claims**

743 Question: Do the main claims made in the abstract and introduction accurately reflect the  
744 paper's contributions and scope?

745 Answer: [Yes]

746 Justification: Claims are supported by both the theoretical results presented in the paper or  
747 experimental exploration.

748 Guidelines:

- 749 • The answer NA means that the abstract and introduction do not include the claims  
750 made in the paper.
- 751 • The abstract and/or introduction should clearly state the claims made, including the  
752 contributions made in the paper and important assumptions and limitations. A No or  
753 NA answer to this question will not be perceived well by the reviewers.
- 754 • The claims made should match theoretical and experimental results, and reflect how  
755 much the results can be expected to generalize to other settings.
- 756 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
757 are not attained by the paper.

758 **2. Limitations**

759 Question: Does the paper discuss the limitations of the work performed by the authors?

760 Answer: [Yes]

761 Justification: Limitations of the proposed embedding method are discussed. For instance,  
762 the requirement of only embedding the MDT are discussed. Impacts on classification due  
763 to this are explored in the “MDT vs DT classification” discussion, with further details in  
764 Appendix. A last paragraph in the conclusion also discusses limitations.

765 Guidelines:

- 766 • The answer NA means that the paper has no limitation while the answer No means that  
767 the paper has limitations, but those are not discussed in the paper.
- 768 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 769 • The paper should point out any strong assumptions and how robust the results are to  
770 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
771 model well-specification, asymptotic approximations only holding locally). The authors  
772 should reflect on how these assumptions might be violated in practice and what the  
773 implications would be.
- 774 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
775 only tested on a few datasets or with a few runs. In general, empirical results often  
776 depend on implicit assumptions, which should be articulated.
- 777 • The authors should reflect on the factors that influence the performance of the approach.  
778 For example, a facial recognition algorithm may perform poorly when image resolution  
779 is low or images are taken in low lighting. Or a speech-to-text system might not be  
780 used reliably to provide closed captions for online lectures because it fails to handle  
781 technical jargon.
- 782 • The authors should discuss the computational efficiency of the proposed algorithms  
783 and how they scale with dataset size.
- 784 • If applicable, the authors should discuss possible limitations of their approach to  
785 address problems of privacy and fairness.
- 786 • While the authors might fear that complete honesty about limitations might be used by  
787 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
788 limitations that aren't acknowledged in the paper. The authors should use their best  
789 judgment and recognize that individual actions in favor of transparency play an impor-  
790 tant role in developing norms that preserve the integrity of the community. Reviewers  
791 will be specifically instructed to not penalize honesty concerning limitations.

### 792 3. Theory Assumptions and Proofs

793 Question: For each theoretical result, does the paper provide the full set of assumptions and  
794 a complete (and correct) proof?

795 Answer: [Yes]

796 Justification: Assumptions are detailed in formal results and full proofs are provided in the  
797 Appendix.

798 Guidelines:

- 799 • The answer NA means that the paper does not include theoretical results.
- 800 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
801 referenced.
- 802 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 803 • The proofs can either appear in the main paper or the supplemental material, but if  
804 they appear in the supplemental material, the authors are encouraged to provide a short  
805 proof sketch to provide intuition.
- 806 • Inversely, any informal proof provided in the core of the paper should be complemented  
807 by formal proofs provided in appendix or supplemental material.
- 808 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 809 4. Experimental Result Reproducibility

810 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
811 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
812 of the paper (regardless of whether the code and data are provided or not)?

813 Answer: [Yes]

814 Justification: Detailed settings of constructing visualizations are presented. Code provided  
815 with an example public domain and a resource file + README.txt for quick testing and  
816 validation.

817 Guidelines:

- 818 • The answer NA means that the paper does not include experiments.
- 819 • If the paper includes experiments, a No answer to this question will not be perceived  
820 well by the reviewers: Making the paper reproducible is important, regardless of  
821 whether the code and data are provided or not.
- 822 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
823 to make their results reproducible or verifiable.
- 824 • Depending on the contribution, reproducibility can be accomplished in various ways.  
825 For example, if the contribution is a novel architecture, describing the architecture fully  
826 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
827 be necessary to either make it possible for others to replicate the model with the same  
828 dataset, or provide access to the model. In general, releasing code and data is often  
829 one good way to accomplish this, but reproducibility can also be provided via detailed  
830 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
831 of a large language model), releasing of a model checkpoint, or other means that are  
832 appropriate to the research performed.
- 833 • While NeurIPS does not require releasing code, the conference does require all submis-  
834 sions to provide some reasonable avenue for reproducibility, which may depend on the  
835 nature of the contribution. For example
  - 836 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
837 to reproduce that algorithm.
  - 838 (b) If the contribution is primarily a new model architecture, the paper should describe  
839 the architecture clearly and fully.
  - 840 (c) If the contribution is a new model (e.g., a large language model), then there should  
841 either be a way to access this model for reproducing the results or a way to reproduce  
842 the model (e.g., with an open-source dataset or instructions for how to construct  
843 the dataset).
  - 844 (d) We recognize that reproducibility may be tricky in some cases, in which case  
845 authors are welcome to describe the particular way they provide for reproducibility.  
846 In the case of closed-source models, it may be that access to the model is limited in  
847 some way (e.g., to registered users), but it should be possible for other researchers  
848 to have some path to reproducing or verifying the results.

## 849 5. Open access to data and code

850 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
851 tions to faithfully reproduce the main experimental results, as described in supplemental  
852 material?

853 Answer: [Yes]

854 Justification: Code provided with an example public domain and a resource file +  
855 README.txt for quick testing and validation.

856 Guidelines:

- 857 • The answer NA means that paper does not include experiments requiring code.
- 858 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
859 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 860 • While we encourage the release of code and data, we understand that this might not be  
861 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
862 including code, unless this is central to the contribution (e.g., for a new open-source  
863 benchmark).
- 864 • The instructions should contain the exact command and environment needed to run to  
865 reproduce the results. See the NeurIPS code and data submission guidelines ([https://nips.cc/  
866 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 867 • The authors should provide instructions on data access and preparation, including how  
868 to access the raw data, preprocessed data, intermediate data, and generated data, etc.

- 869           • The authors should provide scripts to reproduce all experimental results for the new  
870           proposed method and baselines. If only a subset of experiments are reproducible, they  
871           should state which ones are omitted from the script and why.  
872           • At submission time, to preserve anonymity, the authors should release anonymized  
873           versions (if applicable).  
874           • Providing as much information as possible in supplemental material (appended to the  
875           paper) is recommended, but including URLs to data and code is permitted.

876       **6. Experimental Setting/Details**

877       Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
878           parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
879           results?

880       Answer: [Yes]

881       Justification: Details are provided in main text and Appendix.

882       Guidelines:

- 883           • The answer NA means that the paper does not include experiments.  
884           • The experimental setting should be presented in the core of the paper to a level of detail  
885           that is necessary to appreciate the results and make sense of them.  
886           • The full details can be provided either with the code, in appendix, or as supplemental  
887           material.

888       **7. Experiment Statistical Significance**

889       Question: Does the paper report error bars suitably and correctly defined or other appropriate  
890           information about the statistical significance of the experiments?

891       Answer: [Yes]

892       Justification: Error is stated whenever suitable. Multiple results are presented for visualiza-  
893           tions.

894       Guidelines:

- 895           • The answer NA means that the paper does not include experiments.  
896           • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
897           dence intervals, or statistical significance tests, at least for the experiments that support  
898           the main claims of the paper.  
899           • The factors of variability that the error bars are capturing should be clearly stated (for  
900           example, train/test split, initialization, random drawing of some parameter, or overall  
901           run with given experimental conditions).  
902           • The method for calculating the error bars should be explained (closed form formula,  
903           call to a library function, bootstrap, etc.).  
904           • The assumptions made should be given (e.g., Normally distributed errors).  
905           • It should be clear whether the error bar is the standard deviation or the standard error  
906           of the mean.  
907           • It is OK to report 1-sigma error bars, but one should state it. The authors should  
908           preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
909           of Normality of errors is not verified.  
910           • For asymmetric distributions, the authors should be careful not to show in tables or  
911           figures symmetric error bars that would yield results that are out of range (e.g. negative  
912           error rates).  
913           • If error bars are reported in tables or plots, The authors should explain in the text how  
914           they were calculated and reference the corresponding figures or tables in the text.

915       **8. Experiments Compute Resources**

916       Question: For each experiment, does the paper provide sufficient information on the com-  
917           puter resources (type of compute workers, memory, time of execution) needed to reproduce  
918           the experiments?

919       Answer: [NA].

920       Justification: The code runs on any standard computer.

921 Guidelines:

- 922 • The answer NA means that the paper does not include experiments.  
923 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
924 or cloud provider, including relevant memory and storage.  
925 • The paper should provide the amount of compute required for each of the individual  
926 experimental runs as well as estimate the total compute.  
927 • The paper should disclose whether the full research project required more compute  
928 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
929 didn't make it into the paper).

930 **9. Code Of Ethics**

931 Question: Does the research conducted in the paper conform, in every respect, with the  
932 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

933 Answer: [Yes]

934 Justification: The research of the paper follows the code of ethics.

935 Guidelines:

- 936 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.  
937 • If the authors answer No, they should explain the special circumstances that require a  
938 deviation from the Code of Ethics.  
939 • The authors should make sure to preserve anonymity (e.g., if there is a special consider-  
940 ation due to laws or regulations in their jurisdiction).

941 **10. Broader Impacts**

942 Question: Does the paper discuss both potential positive societal impacts and negative  
943 societal impacts of the work performed?

944 Answer: [Yes]

945 Justification: Broader impact and societal impact is briefly mentioned throughout the main  
946 text, with further discussion presented in an Appendix section.

947 Guidelines:

- 948 • The answer NA means that there is no societal impact of the work performed.  
949 • If the authors answer NA or No, they should explain why their work has no societal  
950 impact or why the paper does not address societal impact.  
951 • Examples of negative societal impacts include potential malicious or unintended uses  
952 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
953 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
954 groups), privacy considerations, and security considerations.  
955 • The conference expects that many papers will be foundational research and not tied  
956 to particular applications, let alone deployments. However, if there is a direct path to  
957 any negative applications, the authors should point it out. For example, it is legitimate  
958 to point out that an improvement in the quality of generative models could be used to  
959 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
960 that a generic algorithm for optimizing neural networks could enable people to train  
961 models that generate Deepfakes faster.  
962 • The authors should consider possible harms that could arise when the technology is  
963 being used as intended and functioning correctly, harms that could arise when the  
964 technology is being used as intended but gives incorrect results, and harms following  
965 from (intentional or unintentional) misuse of the technology.  
966 • If there are negative societal impacts, the authors could also discuss possible mitigation  
967 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
968 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
969 feedback over time, improving the efficiency and accessibility of ML).

970 **11. Safeguards**

971 Question: Does the paper describe safeguards that have been put in place for responsible  
972 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
973 image generators, or scraped datasets)?

974                  Answer: [NA]

975                  Justification: No release of data or models.

976                  Guidelines:

- 977                  • The answer NA means that the paper poses no such risks.
- 978                  • Released models that have a high risk for misuse or dual-use should be released with  
979                  necessary safeguards to allow for controlled use of the model, for example by requiring  
980                  that users adhere to usage guidelines or restrictions to access the model or implementing  
981                  safety filters.
- 982                  • Datasets that have been scraped from the Internet could pose safety risks. The authors  
983                  should describe how they avoided releasing unsafe images.
- 984                  • We recognize that providing effective safeguards is challenging, and many papers do  
985                  not require this, but we encourage authors to take this into account and make a best  
986                  faith effort.

## 987                  12. Licenses for existing assets

988                  Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
989                  the paper, properly credited and are the license and terms of use explicitly mentioned and  
990                  properly respected?

991                  Answer: [Yes]

992                  Justification: Appropriate credit / references are provided. Licenses listed in Appendix.

993                  Guidelines:

- 994                  • The answer NA means that the paper does not use existing assets.
- 995                  • The authors should cite the original paper that produced the code package or dataset.
- 996                  • The authors should state which version of the asset is used and, if possible, include a  
997                  URL.
- 998                  • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 999                  • For scraped data from a particular source (e.g., website), the copyright and terms of  
1000                 service of that source should be provided.
- 1001                 • If assets are released, the license, copyright information, and terms of use in the  
1002                 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
1003                 has curated licenses for some datasets. Their licensing guide can help determine the  
1004                 license of a dataset.
- 1005                 • For existing datasets that are re-packaged, both the original license and the license of  
1006                 the derived asset (if it has changed) should be provided.
- 1007                 • If this information is not available online, the authors are encouraged to reach out to  
1008                 the asset's creators.

## 1009                 13. New Assets

1010                 Question: Are new assets introduced in the paper well documented and is the documentation  
1011                 provided alongside the assets?

1012                 Answer: [NA]

1013                 Justification: No new assets provided.

1014                 Guidelines:

- 1015                 • The answer NA means that the paper does not release new assets.
- 1016                 • Researchers should communicate the details of the dataset/code/model as part of their  
1017                 submissions via structured templates. This includes details about training, license,  
1018                 limitations, etc.
- 1019                 • The paper should discuss whether and how consent was obtained from people whose  
1020                 asset is used.
- 1021                 • At submission time, remember to anonymize your assets (if applicable). You can either  
1022                 create an anonymized URL or include an anonymized zip file.

## 1023                 14. Crowdsourcing and Research with Human Subjects

1024 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1025 include the full text of instructions given to participants and screenshots, if applicable, as  
1026 well as details about compensation (if any)?

1027 Answer: [NA]

1028 Justification: No crowdsourcing or research with human subjects.

1029 Guidelines:

- 1030 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1031 human subjects.
- 1032 • Including this information in the supplemental material is fine, but if the main contribu-  
1033 tion of the paper involves human subjects, then as much detail as possible should be  
1034 included in the main paper.
- 1035 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1036 or other labor should be paid at least the minimum wage in the country of the data  
1037 collector.

1038 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
1039 Subjects**

1040 Question: Does the paper describe potential risks incurred by study participants, whether  
1041 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1042 approvals (or an equivalent approval/review based on the requirements of your country or  
1043 institution) were obtained?

1044 Answer: [NA]

1045 Justification: No research with human subjects.

1046 Guidelines:

- 1047 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1048 human subjects.
- 1049 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1050 may be required for any human subjects research. If you obtained IRB approval, you  
1051 should clearly state this in the paper.
- 1052 • We recognize that the procedures for this may vary significantly between institutions  
1053 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1054 guidelines for their institution.
- 1055 • For initial submissions, do not include any information that would break anonymity (if  
1056 applicable), such as the institution conducting the review.