



# Interpretable Image Classification Using Sparse Oblique Decision Trees



Suryabhan Singh Hada and Miguel Á. Carreira-Perpiñán,  
Dept. CSE, UC Merced

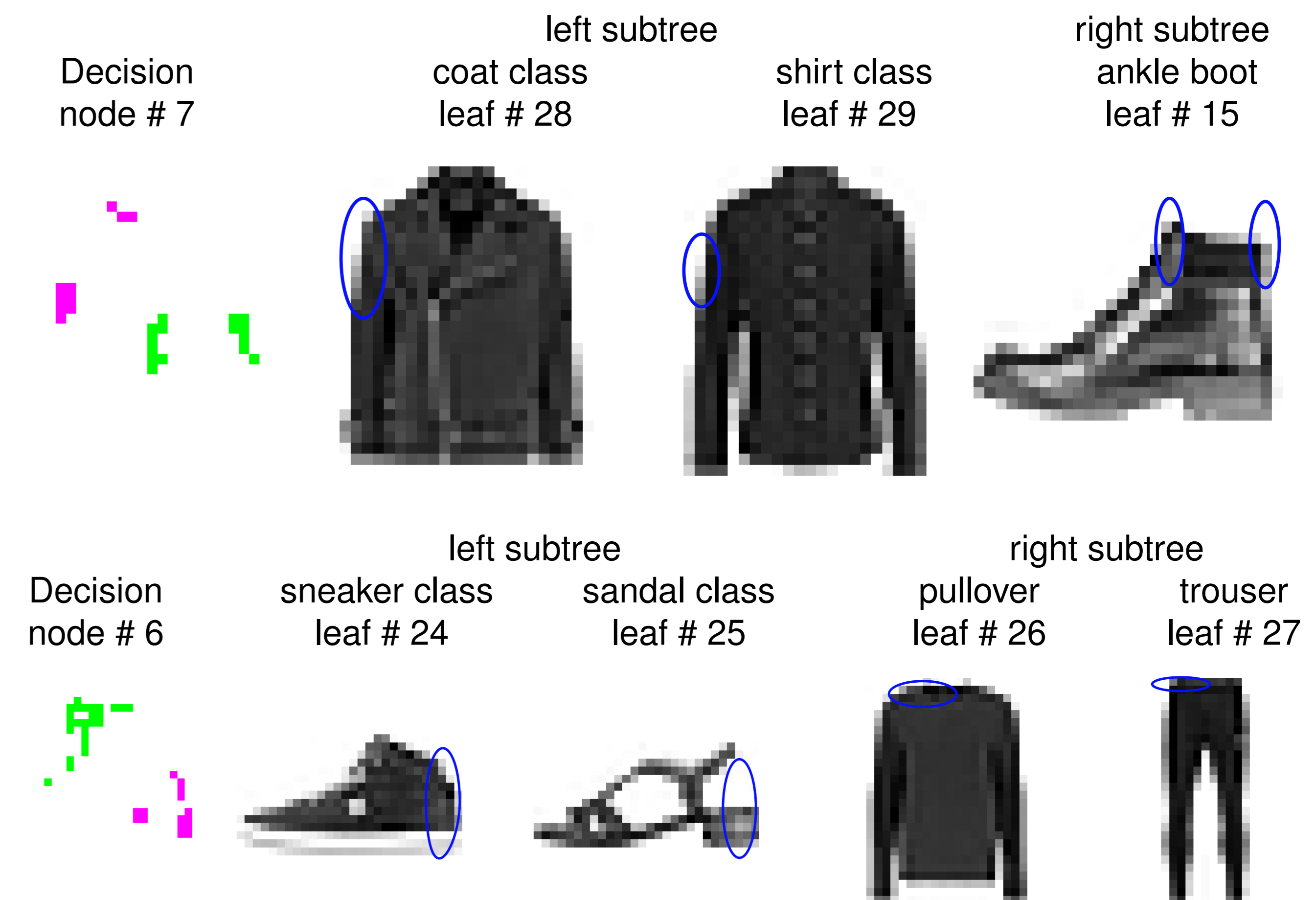
## 1 Motivation and summary

- Interpreting the image datasets is a difficult task, as each image contains a lot of irrelevant data. This makes it hard to understand what part of the image is important or what common concept defines a particular category of the class.
- We address these issues by using sparse oblique trees as a tool to understand the given image dataset. Unlike axis-aligned trees that operate only on a single feature at each node, the sparse oblique tree operates on a small, learnable subset of features.
- Sparse oblique trees are not only accurate but also very interpretable. We can learn accurate enough sparse oblique trees with the tree alternating optimization (TAO) algorithm.

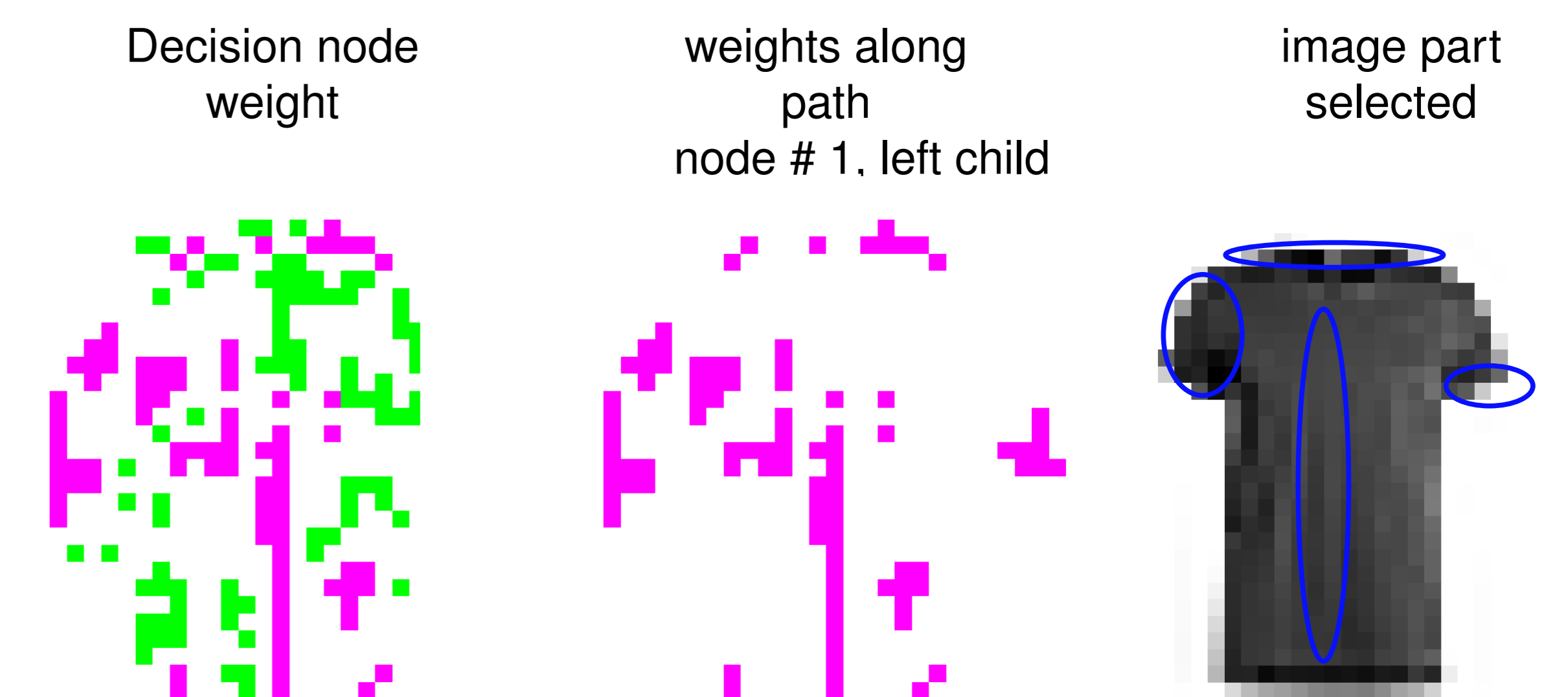
## 2 Interpret dataset using sparse oblique tree

- Train a sparse oblique tree using TAO and pick the sparsity parameter such that the resultant tree is as sparse as possible but remains accurate enough. Next, inspect the weights ( $\mathbf{w}$ ) of the decision nodes to extract relevant features from the dataset.
- Write weight vector  $\mathbf{w}$  and input  $\mathbf{x}$  as  $\mathbf{w} = (\mathbf{w}_0 \ \mathbf{w}_- \ \mathbf{w}_+)$  and  $\mathbf{x} = (\mathbf{x}_0 \ \mathbf{x}_- \ \mathbf{x}_+)$ , where  $\mathbf{w}_0 = \mathbf{0}$ ,  $\mathbf{w}_- < \mathbf{0}$  and  $\mathbf{w}_+ > \mathbf{0}$  contain the zero, negative and positive weights in  $\mathbf{w}$ , and  $\mathbf{x}$  is arranged accordingly. Call  $S_0$ ,  $S_-$  and  $S_+$  the corresponding sets of indices in  $\mathbf{w}$ .
  - ◊ If  $\mathbf{x}$  goes right, we represent the feature selected as a binary vector  $\mu_+ \in \{0, 1\}^d$ , containing ones only at  $S_+$ .
  - ◊ If  $\mathbf{x}$  goes left, we represent the feature selected as a binary vector  $\mu_- \in \{0, 1\}^d$ , containing ones only at  $S_-$ .
- We call  $\mu^+$  and  $\mu^-$  the NODE-FEATURES, where location of one represents features selected by  $\mathbf{w}$ . We use NODE-FEATURES, to interpret the dataset as follows:
  - ◊ For each decision node NODE-FEATURES represents the features related to left and right subtree. By using NODE-FEATURES, we can understand what set of features separate a group of classes.
  - ◊ Features associated with a class  $k$ : for each node in the path from the root to leaf for class  $k$  collect NODE-FEATURES, and at the end take logical OR of all NODE-FEATURES. If there is more than one leaf for class  $k$ , take the union of all the features selected.
  - ◊ For features specific to a given an input  $\mathbf{x}$ : Repeat the process as above, but only for the leaf containing the input  $\mathbf{x}$ . Next, keep only those features that are active in the  $\mathbf{x}$ .
- We can plot these features to visualize what concept is captured by these features.

Difference among group of classes



Feature selection for a given instance



Fashion-MNIST dataset

