

# Regularizing Soft Decision Trees

Olcay Taner Yıldız<sup>1</sup> Ethem Alpaydın<sup>2</sup>

<sup>1</sup>Dept of Computer Engineering, Işık University, TR-34980, Istanbul, Turkey

<sup>2</sup>Dept of Computer Engineering, Boğaziçi University, TR-34342, Istanbul, Turkey

ISCIS 2013

# Outline

- 1 Introduction
- 2 Soft Decision Tree
- 3 Regularization
- 4 Experiments
- 5 Conclusions

# Hard Decision Tree

- Each decision node  $m$  applies a test  $g_m(\mathbf{x})$  and chooses one of the children accordingly.

$$F_m(\mathbf{x}) = \begin{cases} F_m^L(\mathbf{x}) & \text{if } g_m(\mathbf{x}) > 0 \text{ /* true */} \\ F_m^R(\mathbf{x}) & \text{otherwise /* false */} \end{cases}$$

- Classification: Leaves carry the label of one of  $K$  classes
- Regression: Leaves carry a constant which is the numeric regression value.

## Hard Decision Tree Types

- *Univariate tree*:  $g_m(\mathbf{x}) = x_j + w_{m0} > 0$ .
- *Multivariate linear tree*:  $g_m(\mathbf{x}) = \mathbf{w}_m^T \mathbf{x} + w_{m0} > 0$ .
- *Multivariate nonlinear tree*:  $g_m(\mathbf{x}) = \sum_{j=1}^k w_j \phi_j(\mathbf{x}) > 0$ .
- *Omnivariate tree*:  $g_m(\mathbf{x})$  can be any of the above, chosen by a statistical model selection procedure.

# Soft Decision Tree

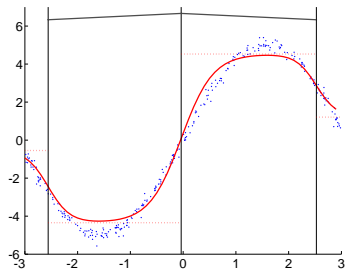
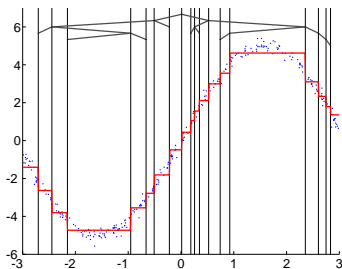
- Soft decision node redirects instances to all its children with probabilities calculated by a *gating function*  $g_m(\mathbf{x})$ .

$$F_m(\mathbf{x}) = F_m^L(\mathbf{x})g_m(\mathbf{x}) + F_m^R(\mathbf{x})(1 - g_m(\mathbf{x}))$$

$$g_m(\mathbf{x}) = \frac{1}{1 + \exp[-(\mathbf{w}_m^T \mathbf{x} + w_{m0})]}$$

- Gating model implements a discriminative (logistic linear) model estimating the posterior probability of the left child.

# Hard vs. Soft Tree (Toy Dataset)



## Response and Error

```
1 function  $F_m(\mathbf{x})$ 
2   if  $m$  is leaf node
3      $y = z_m$  /* leaf value at  $m$  */
4   else
5      $g_m(\mathbf{x}) = \text{sigmoid}(\mathbf{w}_m^T \mathbf{x} + w_{m0})$ 
6      $y = F_m^L(\mathbf{x})g_m(\mathbf{x}) + F_m^R(\mathbf{x})(1 - g_m(\mathbf{x}))$ 
7   return  $y$ 
```

- Classification:  $E = r \log y + (1 - r) \log(1 - y)$

# Training the Soft Decision Tree

```

1 function LearnSoftTree( $m, \mathcal{X}, \mathcal{V}$ )
2    $E_{before} = \text{ErrorOfTree}(\mathcal{V})$ 
3   initialize  $w_{mj}, z_m^L$ , and  $z_m^R$ 
4   repeat
5     for all  $(\mathbf{x}, r) \in \mathcal{X}$ 
6        $\delta(\mathbf{x}) = (F_{root}(\mathbf{x}) - r)(g_p(\mathbf{x}))^{left}(1 - g_p(\mathbf{x}))^{right}$ 
7       for  $j = 0, \dots, d$ 
8          $w_{mj} = w_{mj} - \eta \delta(\mathbf{x})(F_m^L(\mathbf{x}) - F_m^R(\mathbf{x}))v_m(\mathbf{x})(1 - v_m(\mathbf{x}))x_j$ 
9          $z_m^L = z_m^L - \eta \delta(\mathbf{x})v_m(\mathbf{x})$ 
10         $z_m^R = z_m^R - \eta \delta(\mathbf{x})(1 - v_m(\mathbf{x}))$ 
11    until convergence
12     $E_{after} = \text{ErrorOfTree}(\mathcal{V})$ 
13    if  $E_{after} < E_{before}$ 
14      LearnSoftTree( $m.left, \mathcal{X}, \mathcal{V}$ )
15      LearnSoftTree( $m.right, \mathcal{X}, \mathcal{V}$ )
  
```



# Regularization

- Local dimensionality reduction for better generalization
- $L_1$  regularization:

$$E_{L_1} = (1 - \lambda) \text{CrossEntropy} + \lambda \sum_{i=0}^d |w_{mi}|$$

- $L_2$  regularization:

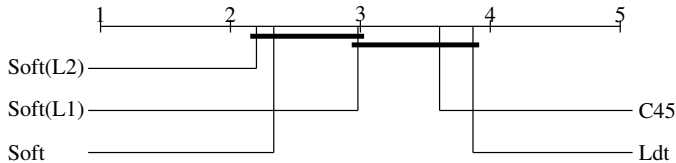
$$E_{L_2} = (1 - \lambda) \text{CrossEntropy} + \lambda \sum_{i=0}^d w_{mi}^2$$

# Experiments

- Hard decision tree (HDT), Linear discriminant tree (LDT), Soft decision tree (SDT) without regularization, SDT with  $L_1$  regularization, SDT with  $L_2$  regularization.
- Comparison on 27 data sets
- 2/3 training data with  $5 \times 2$ -fold cross validation, 1/3 test set.
- Parametric  $5 \times 2$  paired  $F$ -test used for comparison on a single data set and nonparametric Nemenyi's test for overall comparison.

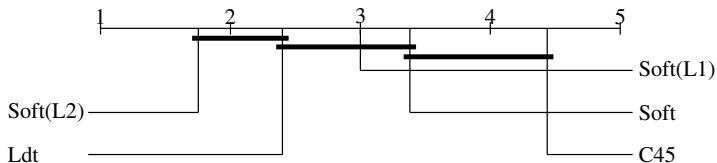
# Results: Accuracy

	Hard	Ldt	Soft	Soft( $L_1$ )	Soft( $L_2$ )
Hard		6	0	5	4
Ldt	5		0	1	1
Soft	9	11		5	4
Soft( $L_1$ )	7	4	1		0
Soft( $L_2$ )	8	4	1	2	



## Results: Node Counts

	Hard	Ldt	Soft	Soft( $L_1$ )	Soft( $L_2$ )
Hard		0	1	1	1
Ldt	11		4	2	0
Soft	10	2		1	0
Soft( $L_1$ )	12	3	2		0
Soft( $L_2$ )	13	4	4	1	



## Conclusions: Soft Trees

- Proposed decision tree model with soft decisions, which makes use of a soft gating function to merge the decisions of the subtrees.
- Soft trees have smoother fits and hence lower bias around the split boundaries.
- Linear gating function enables soft trees to make oblique splits in contrast to the axis-orthogonal splits made by hard trees.

## Conclusions: Regularization

- We extend the soft decision tree model by adding  $L_1$  and  $L_2$  regularization to penalize unnecessary complexity.
- Both versions improve accuracy slightly and decrease complexity significantly.
- $L_2$  regularization seems to work slightly better than  $L_1$ .