# Efficient Construction of Decision Trees by the Dual Information Distance Method

Irad Ben-Gal
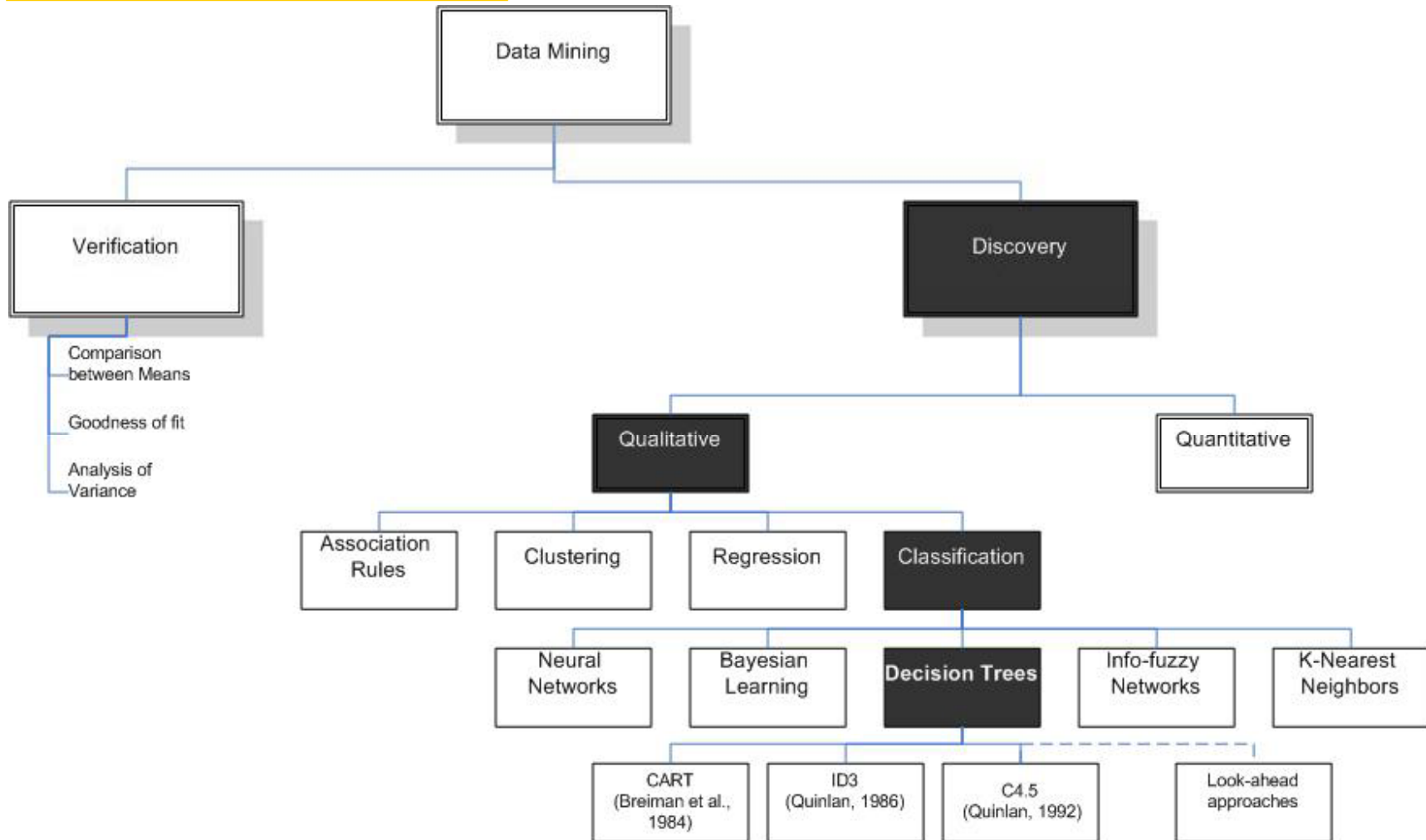
Joint work with Alexandra Dana, Niv Shkolnik and Gonen Singer

Dept. of Industrial Engineering

Tel Aviv University

# Outline

# Introduction

Data Mining

Verification
- Comparison between Means
- Goodness of fit
- Analysis of Variance

Discovery

Qualitative

Quantitative

Association Rules

Clustering

Regression

Classification

Neural Networks

Bayesian Learning

**Decision Trees**

Info-fuzzy Networks

K-Nearest Neighbors

CART (Breiman et al., 1984)

ID3 (Quinlan, 1986)

C4.5 (Quinlan, 1992)

Look-ahead approaches

Based on Maimon & Rokach (2005)

# Introduction

Classification (Supervised Learning)

"In classification, there is a target categorical variable, which is partitioned into predetermined classes or categories. The data mining model examine a large set of records, each record containing information on the target variable as well as a set of input or predictor variables". (Larose, 2005).
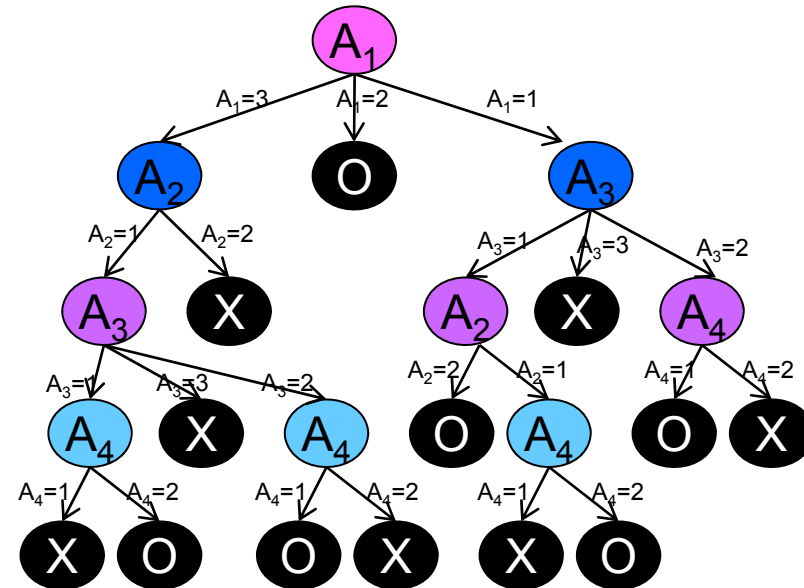
# Introduction

> **Construction of Decision Tree**

❑ Classification variables

- Class/Target variable Y
- Attributes set $\{A_1, A_2, …, A_n\}$

❑ The tree partitions Y by selecting attributes $A_i$, aiming that each leaf will contain a single class of Y

❑ Performance measures

- Minimizing the classification error-rates
- Minimizing the average depth of the tree
- Minimizing the number of nodes/leaves

# Introduction

| #  | A1 | A2 | A3 | A4 | Y |
|----|----|----|----|----|---|
| 1  | 1  | 1  | 1  | 1  | X |
| 2  | 1  | 1  | 1  | 2  | O |
| 3  | 1  | 1  | 2  | 1  | O |
| 4  | 1  | 1  | 2  | 2  | X |
| 5  | 1  | 1  | 3  | 2  | X |
| 6  | 1  | 2  | 1  | 1  | O |
| 7  | 1  | 2  | 1  | 2  | O |
| 8  | 1  | 2  | 3  | 2  | X |
| 9  | 2  | 1  | 1  | 1  | O |
| 10 | 2  | 1  | 1  | 2  | O |
| 11 | 2  | 1  | 2  | 1  | O |
| 12 | 2  | 2  | 1  | 1  | O |
| 13 | 2  | 2  | 1  | 2  | O |
| 14 | 2  | 2  | 2  | 1  | O |
| 15 | 2  | 2  | 2  | 2  | O |
| 16 | 3  | 1  | 1  | 2  | X |
| 17 | 3  | 1  | 2  | 1  | X |
| 18 | 3  | 1  | 2  | 2  | O |
| 19 | 3  | 1  | 3  | 1  | X |
| 20 | 3  | 1  | 3  | 2  | X |
| 21 | 3  | 2  | 2  | 1  | X |
| 22 | 3  | 2  | 2  | 2  | X |
| 23 | 3  | 2  | 3  | 1  | X |
| 24 | 3  | 1  | 1  | 1  | O |
| 25 | 3  | 2  | 3  | 2  | X |

## Example: ID3 (Lee & Olafsson, 2006)



ID3 Decision Tree example

- 4 categorical attributes $A_1, \ldots, A_4$
- Each selection splits the set into two or more partitions

# Introduction

| | A1 | A2 | A3 | A4 | Y |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 | X |
| 2 | 1 | 1 | 1 | 2 | O |
| 3 | 1 | 1 | 2 | 1 | O |
| 4 | 1 | 1 | 2 | 2 | X |
| 5 | 1 | 1 | 3 | 2 | X |
| 6 | 1 | 2 | 1 | 1 | O |
| 7 | 1 | 2 | 1 | 2 | O |
| 8 | 1 | 2 | 3 | 2 | X |
| 9 | 2 | 1 | 1 | 1 | O |
| 10 | 2 | 1 | 1 | 2 | O |
| 11 | 2 | 1 | 2 | 1 | O |
| 12 | 2 | 2 | 1 | 1 | O |
| 13 | 2 | 2 | 1 | 2 | O |
| 14 | 2 | 2 | 2 | 1 | O |
| 15 | 2 | 2 | 2 | 2 | O |
| 16 | 3 | 1 | 1 | 2 | X |
| 17 | 3 | 1 | 2 | 1 | X |
| 18 | 3 | 1 | 2 | 2 | O |
| 19 | 3 | 1 | 3 | 1 | X |
| 20 | 3 | 1 | 3 | 2 | X |
| 21 | 3 | 2 | 2 | 1 | X |
| 22 | 3 | 2 | 2 | 2 | X |
| 23 | 3 | 2 | 3 | 1 | X |
| 24 | 3 | 1 | 1 | 1 | O |
| 25 | 3 | 2 | 3 | 2 | X |

**Partition by Y**

$$\alpha_Y \equiv \left\langle \begin{array}{l} \{1,4,5,8,16,17,19,20,21,22,23,25\}, \\ \{2,3,6,7,9,10,11,12,13,14,15,18,24\} \end{array} \right\rangle$$

# Introduction

| | A1 | A2 | A3 | A4 | Y |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 | X |
| 2 | 1 | 1 | 1 | 2 | O |
| 3 | 1 | 1 | 2 | 1 | O |
| 4 | 1 | 1 | 2 | 2 | X |
| 5 | 1 | 1 | 3 | 2 | X |
| 6 | 1 | 2 | 1 | 1 | O |
| 7 | 1 | 2 | 1 | 2 | O |
| 8 | 1 | 2 | 3 | 2 | X |
| 9 | 2 | 1 | 1 | 1 | O |
| 10 | 2 | 1 | 1 | 2 | O |
| 11 | 2 | 1 | 2 | 1 | O |
| 12 | 2 | 2 | 1 | 1 | O |
| 13 | 2 | 2 | 1 | 2 | O |
| 14 | 2 | 2 | 2 | 1 | O |
| 15 | 2 | 2 | 2 | 2 | O |
| 16 | 3 | 1 | 1 | 2 | X |
| 17 | 3 | 1 | 2 | 1 | X |
| 18 | 3 | 1 | 2 | 2 | O |
| 19 | 3 | 1 | 3 | 1 | X |
| 20 | 3 | 1 | 3 | 2 | X |
| 21 | 3 | 2 | 2 | 1 | X |
| 22 | 3 | 2 | 2 | 2 | X |
| 23 | 3 | 2 | 3 | 1 | X |
| 24 | 3 | 1 | 1 | 1 | O |
| 25 | 3 | 2 | 3 | 2 | X |

**Partition by** $A_1$

$$\alpha_1 = \left\langle \begin{array}{l} \{1,2,3,4,5,6,7,8\}, \\ \{9,10,11,12,13,14,15\}, \\ \{16,17,18,19,20,21,22,23,24,25\} \end{array} \right\rangle$$

# Introduction

| | A1 | A2 | A3 | A4 | Y |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 | X |
| 2 | 1 | 1 | 1 | 2 | O |
| 3 | 1 | 1 | 2 | 1 | O |
| 4 | 1 | 1 | 2 | 2 | X |
| 5 | 1 | 1 | 3 | 2 | X |
| 6 | 1 | 2 | 1 | 1 | O |
| 7 | 1 | 2 | 1 | 2 | O |
| 8 | 1 | 2 | 3 | 2 | X |
| 9 | 2 | 1 | 1 | 1 | O |
| 10 | 2 | 1 | 1 | 2 | O |
| 11 | 2 | 1 | 2 | 1 | O |
| 12 | 2 | 2 | 1 | 1 | O |
| 13 | 2 | 2 | 1 | 2 | O |
| 14 | 2 | 2 | 2 | 1 | O |
| 15 | 2 | 2 | 2 | 2 | O |
| 16 | 3 | 1 | 1 | 2 | X |
| 17 | 3 | 1 | 2 | 1 | X |
| 18 | 3 | 1 | 2 | 2 | O |
| 19 | 3 | 1 | 3 | 1 | X |
| 20 | 3 | 1 | 3 | 2 | X |
| 21 | 3 | 2 | 2 | 1 | X |
| 22 | 3 | 2 | 2 | 2 | X |
| 23 | 3 | 2 | 3 | 1 | X |
| 24 | 3 | 1 | 1 | 1 | O |
| 25 | 3 | 2 | 3 | 2 | X |

**Partition by** $A_1 \vee A_2$

$$a_{12} = \left\langle \begin{array}{l} \{1,2,3,4,5\}, \{6,7,8\}, \\ \{9,10,11\}, \{12,13,14,15\}, \\ \{16,17,18,19,20,24\}, \{21,22,23,25\} \end{array} \right\rangle$$

# Introduction

## Optimal Decision Trees

❑ Consider all possible splits (all combinations)

❑ Construction is NP-hard (Hackock et al., 1996)

## Heuristic Trees

❑ Greedy trees (ID3, C4.5, CART)

- At each step consider only the next split

❑ Look-ahead trees

- consider up-coming splits (usually 2-steps ahead)
- Computationally "pricy": $O(mn^K)$ for n variables; m records, and a K-steps look-ahead procedure

# Introduction

## Types of Decision Tree Algorithms

### Heuristic Trees

Can we do better than greedy selection without extensive computations?

Maybe: e.g. by using the Dual Information Distance (DID) approach on a 'partitions graph'

# Introduction

ID3 **(Quinlan, 1986) &** C4.5 **(Quinlan, 1993)**

❑ Recursively split each node until no splits are possible
❑ ID3 Splitting Criteria: (highest) Information Gain

$$\text{Information Gain}(Y; A_i) = H(Y) - H(Y \mid A_i)$$

❑ C4.5 Splitting Criteria: (highest) Gain Ratio

$$\text{Gain Ratio}(Y; A_i) = \frac{\text{Information Gain}(Y; A_i)}{\text{Enropy}(A_i)}$$

❑ Max information gain/Reduction of entropy (uncertainty)!
❑ Very popular & produces good results in practice problems (Goodman & Smyth, 1988; Murthy & Salzberg, 1995)

# Outline

1. Introduction & Motivation

2. **IT Partitions Approach**

3. Example

4. Results

5. Mid-level solutions

6. Summary & Contribution

# Information Theory

For two attributes, $A_1, A_2$ :

$$I(A_1, A_2; Y) = I(A_1; Y) + I(A_2; Y/A_1) =$$

$$H(A_1) + H(A_2/A_1) - H(A_1/Y) - H(A_2/A_1, Y)$$

In General :

$$I(A_1, A_2, ..., A_n; Y) = \sum_i I(A_i; Y/A_{i-1}, ..., A_1) =$$

$$\sum_i \underbrace{H(A_i/A_{i-1}, ..., A_1)}_{\text{Max Orthogonality (DOE)}} - \sum_i \underbrace{H(A_i/A_{i-1}, ..., A_1, Y)}_{\substack{\text{Minimize Entropy} \ / \\ \text{Remaining uncertainty (DM)}}}$$

# Information Theory

**Chain Rule of Mutual Information (Shannon, 48)**

In General :

$$I(A_1, A_2, ..., A_n; Y) = \sum_i I(A_i; Y \mid A_{i-1}, ..., A_1) =$$

$$\sum_i \underbrace{H(A_i \mid A_{i-1}, ..., A_1)}_{\text{Max Orthogonality (DOE)}} - \sum_i \underbrace{H(A_i \mid A_{i-1}, ..., A_1, Y)}_{\substack{\text{Minimize Entropy /} \\ \text{Remaining uncertainty (DM)}}}$$

Use a Dual Information Distance approach, wrt:

1. Current state / partition $\quad \alpha_c = A_{i-1} \vee A_{i-2}, ..., A_1$

2. Target variable

# The DID Partition Approach

**Orthogonality and Information Gain**

- ❑ When selecting attribute $A_i$ examine the resulting partition with respect to dual Inf. distance:
- ✓ Current partition (current tree state)
- ✓ Target/class partition



Maximize — Minimize

$\alpha_C$   $\alpha_i$   $\alpha_Y$

**Orthogonality**
**(DOE approach)**

**Remaining Uncertainty**
**(DM approach)**

# LRTA* Algorithm

$dist_1(s_0, s_1)$

$dist_1(s_0, s_2)$

$dist_1(s_0, s_n)$

$dist_2(s_1, s_Y)$

$dist_2(s_2, s_Y)$

$dist_2(s_n, s_Y)$

$s_0$  $s_1$  $s_2$  $s_3$  $s_n$  $s_{1,2}$  $s_{2,3}$  $s_{3,1}$  $s_{i,j,k}$  $s_Y$

(Korf, 1990)

**Applying LRTA* concepts to Decision-Tree Construction**

**State $s_i$ ↔ partition $\alpha_i$**
**Neighbors of state $s_i$ ↔ Refinement of partition $\alpha_i$ (add an attribute)**
**Distance $d(s_i, s_j)$ ↔ A distance metric defined over partitions**

# The DID Approach

$$Rokhlin(\alpha, \beta) = H(\alpha|\beta) + H(\beta|\alpha)$$

❑ **This distance follows the required properties of a metric:**

$$d(\alpha, \beta) \geq 0$$

$$d(\alpha, \alpha) = 0$$

$$d(\alpha, \beta) \leq d(\alpha, \gamma) + d(\gamma, \beta), \forall \alpha, \beta, \gamma \in \chi$$

❑ **Relation between mutual information and Rokhlin metric**

$$d(\alpha, \beta) = H(\alpha, \beta) - I(\alpha, \beta)$$

## The DID Approach

A **general objective function**

$$\min_{\alpha_j}\{w_1 d_1(\alpha_c, \alpha_j) + w_2 d_2(\alpha_j, \alpha_Y)\}$$

$d_1$   denotes the orthogonality measure distance between the current partition and the next chosen partition

$d_2$   denotes the information (or remaining uncertainty) distance between the chosen partition and the class partition



Maximize     $\alpha_c$     $\alpha_j$     Minimize     $\alpha_Y$

$w_1 < 0$

$w_2 > 0$

Orthogonality
(DOE approach)

Remaining Uncertainty
(DM approach)

Korf, 1990

A **general objective function**

$$\min_{\alpha_j}\{w_1 d_1(\alpha_c, \alpha_j) + w_2 d_2(\alpha_j, \alpha_Y)$$

$d_1$    denotes the orthogonality measure distance between the current partition and the next chosen partition

$d_2$    denotes the information (or remaining uncertainty) distance between the chosen partition and the class partition



$w_1 < 0$

$w_2 > 0$

Rokhlin ->
Conditional Entropy

Rokhlin

Korf, 1990

# The DID Approach

❑ Let us consider the following example:

❑ For $A_1$ partition $A_3$ is most orthogonal

$$H(A_3 \mid A_1) = 0.63 \quad H(A_2 \mid A_1) = 1.24$$

| $A_1$ | $A_2$ | $A_3$ | Y |
|-------|-------|-------|---|
| 1 | 1 | 1 | 2 |
| 1 | 1 | 2 | 2 |
| 1 | 1 | 2 | 2 |
| 1 | 2 | 2 | 3 |
| 1 | 2 | 2 | 3 |
| 1 | 2 | 2 | 3 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 1 |

# The DID Approach

- Let us consider the following example:
- $A_1$ and $A_3$ are most orthogonal, however they **do not classify Y**

| $A_1$ | $A_2$ | $A_3$ | Y |
|-------|-------|-------|---|
| 1 | 1 | 1 | 2 |
| 1 | 1 | 2 | **2** |
| 1 | 1 | 2 | **2** |
| 1 | 2 | 2 | **3** |
| 1 | 2 | 2 | **3** |
| 1 | 2 | 2 | **3** |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 1 |

# The DID Approach

**Orthogonality between restricted partitions**

- Let us consider the following example:
- $A_1$ and $A_3$ are most orthogonal, however they do not classify $Y$
- Focusing on the orthogonality between the restricted partitions:

  For $A_1 =1$: $A_2$ is most orthogonal

  For $A_1 =2$: $A_3$ is most orthogonal but this sub-set is already classified by $A_1$

| $A_1$ | $A_2$ | $A_3$ | $Y$ |
|---|---|---|---|
| 1 | 1 | 1 | 2 |
| 1 | 1 | 2 | 2 |
| 1 | 1 | 2 | 2 |
| 1 | 2 | 2 | 3 |
| 1 | 2 | 2 | 3 |
| 1 | 2 | 2 | 3 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 1 |

# The DID Approach

**Orthogonality between restricted partitions**

- Let us consider the following example:
- $A_1$ and $A_3$ are most orthogonal, however they do not classify $Y$
- Focusing on the orthogonality between the restricted partitions:

  For $A_1 = 1$: $A_2$ is most orthogonal

  For $A_1 = 2$: $A_3$ is most orthogonal but this sub-set is already classified by $A_1$
- Full classification of $Y$ is achieved
- The DM "Preprocess Approach" is not always helpful!

| $A_1$ | $A_2$ | $A_3$ | Y |
|-------|-------|-------|---|
| 1 | 1 | 1 | 2 |
| 1 | 1 | 2 | 2 |
| 1 | 1 | 2 | 2 |
| 1 | 2 | 2 | 3 |
| 1 | 2 | 2 | 3 |
| 1 | 2 | 2 | 3 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 1 | 1 |
| 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 1 |

# The DID Approach

❑ The Rokhlin distance measure takes into account two terms:

$$H(Y|A_i) \qquad H(A_i|Y)$$

| A$_1$ | A$_2$ | Y |
|---|---|---|
| 1 | 1 | X |
| 2 | 1 | X |
| 3 | 1 | X |
| 4 | 2 | O |
| 5 | 2 | O |
| 6 | 2 | O |

# The DID Approach

❑ The Rokhlin distance measure takes into account two terms:

$$H(Y|A_i) \qquad H(A_i|Y)$$

$$H(Y|A_1) = 0$$
$$H(Y|A_2) = 0$$

| $A_1$ | $A_2$ | Y |
|-------|-------|---|
| 1 | 1 | X |
| 2 | 1 | X |
| 3 | 1 | X |
| 4 | 2 | O |
| 5 | 2 | O |
| 6 | 2 | O |

# The DID Approach

❑ The Rokhlin distance measure takes into account two terms:

$$H(Y|A_i) \qquad H(A_i|Y)$$

$$H(Y|A_1) = 0$$
$$H(Y|A_2) = 0$$

$$H(A_1|Y) = 1.58$$
$$H(A_2|Y) = 0$$

| $A_1$ | $A_2$ | Y |
|-------|-------|---|
| 1 | 1 | X |
| 2 | 1 | X |
| 3 | 1 | X |
| 4 | 2 | O |
| 5 | 2 | O |
| 6 | 2 | O |

$$\{\alpha_1^1, \alpha_1^2, ..., \alpha_1^6\},$$

$$\{\alpha_2^1, \alpha_2^2\}$$

❑ Prefer $A_2$ over $A_1$

❑ Partition the dataset "as little" as required

## The DID Approach

**Notation**

- $\alpha_0$ - The "initial partition": $\{\{X\},\{\varnothing\}\}$

- $\alpha_Y$ - The "class/target partition"

- $\alpha_c$ - current partition (current state of the tree)

- $\alpha_i$ - the partition that results from the selection of attribute $i$

- $\alpha_i^j$ - sub-partition of $\alpha_i$ where the levels of attribute $i$ are equal

  to $j$ $\quad \alpha_i = \{\alpha_i^j | \alpha_i^j \in \alpha_i, \alpha_i^j \cap \alpha_i^k = \varnothing, \bigcup_j \alpha_i^j = \alpha_i\}$

- Selecting an attribute in a node may results in a refinement of the current

  partition, i.e., $a_i \vee a_j$ is a refinement of $a_i$

  $\alpha_i \vee \alpha_j = \{\alpha_i^l \cap \alpha_j^k | l = 1,2,\dots,|\alpha_i|, k = 1,2,\dots,|\alpha_j|\}$

# The DID Approach

**Information and Remaining Uncertainty**

❑ Look for a partition that results in maximum information (minimum classification uncertainty)

❑ Using **Entropy**: choosing the partition i which gives minimum $H(a_Y | a_i)$

❑ Using **Rokhlin**: choosing the partition i which gives minimum $R(a_Y, a_i) = H(a_Y | a_i) + H(a_i | a_Y)$

**Minimizing the classification uncertainty**

**staying "as close as possible" to $\alpha_Y$ avoiding unnecessary refinement**

# The Proposed DID Algorithm

**Algorithm (DID)**

Given  i) set of weights $w_1, w_2$; ii) two distance metrics denoted by $d_1$ and $d_2$;

  iii) *attributes partitions* $\alpha_1, \alpha_2, \ldots, \alpha_n$; and iv) a class partition $\alpha_Y$

Do:

 *Init current partition* $\alpha_c \leftarrow \alpha_0$

*Init* $E = \{\emptyset\}, F = \{1, 2, \ldots, n\}$ *$<$groups of the "used" and the "unused" attributes $>$*

  *For each sub-partition* $\alpha_c^i \epsilon \alpha_c$ *such that i)* $\left| \alpha_c^i \right| > 1$; *ii)* $\alpha_Y|_{\alpha_c^i}$ *is not yet*

  *classified; and iii)* $F$ *is not empty*

  *start the* **Search** *procedure (for the sub-partition* $\alpha_c^i$)

# The Proposed DID Algorithm

**Algorithm (DID)**

Given i) set of weights $w_1, w_2$; ii) two distance metrics denoted by $d_1$ and $d_2$;

iii) *attributes partitions* $\alpha_1, \alpha_2, \ldots, \alpha_n$; and iv) a class partition $\alpha_Y$

Do:

*Init current partition* $\alpha_c \leftarrow \alpha_0$

*Init* $E = \{\emptyset\}, F = \{1, 2, \ldots, n\}$ *<groups of the "used" and the "unused" attributes >*

*For each sub-partition* $\alpha_c^i \in \alpha_c$ *such that i)* $\left|\alpha_c^i\right| > 1$; *ii)* $\alpha_Y|_{\alpha_c^i}$ *is not yet*

*classified; and iii) F is not empty*

*start the* **Search** *procedure (for the sub-partition* $\alpha_c^i$)

# The Proposed DID Algorithm

**Algorithm (DID)**

Given i) set of weights $w_1, w_2$; ii) two distance metrics denoted by $d_1$ and $d_2$;

iii) *attributes partitions* $\alpha_1, \alpha_2, \ldots, \alpha_n$; and iv) a class partition $\alpha_Y$

Do:

*Init current partition* $\alpha_c \leftarrow \alpha_0$

*Init* $E = \{\emptyset\}, F = \{1, 2, \ldots, n\}$ *<groups of the "used" and the "unused" attributes >*

*For each sub-partition* $\alpha_c^i \epsilon \alpha_c$ *such that i)* $\left| \alpha_c^i \right| > 1$; *ii)* $\alpha_Y |_{\alpha_c^i}$ *is not yet*

*classified; and iii)* $F$ *is not empty*

*start the* **Search** *procedure (for the sub-partition* $\alpha_c^i$*)*

# The Proposed DID Algorithm

**Function Search**. *Given set $\alpha_c^i$ and the attributes partitions $\alpha_1, \alpha_2, ..., \alpha_n$; $E$; $F$;*

1. *Init current partition $\alpha_c$ by $\alpha_c^i$; init $E_c \leftarrow E$; $F_c \leftarrow F$*

2. *Normalize probabilities of the elements of $\alpha_c^i$.*

3. *Create local class partition $\alpha_Y|_{\alpha_c}$.*

4. *Generate neighborhood partitions:*

    $N(\alpha_c) = \{\alpha_j|_{\alpha_c}\}, j\epsilon F_c, F_c = \{j : A_j$ not selected by the algorithm yet$\}$

5. *Normalize probabilities of neighbors and of the class partition;*

6. *Obtain distance measures by $d_1(\alpha_c, \alpha_j|_{\alpha_c})$ and $d_2(\alpha_j|_{\alpha_c}, \alpha_Y|_{\alpha_c})$, $j\epsilon F_c$*

7. *Choose next partition:* The next partition is selected as follows (ties are resolved arbitrary):

    $\alpha_{next} \leftarrow \arg\min_{\alpha_j \in N(\alpha_c)}\{w_1 d_1(\alpha_c, \alpha_j|_{\alpha_c}) + w_2 d_2(\alpha_j|_{\alpha_c}, \alpha_Y|_{\alpha_c})\}$

8. *Update $E_c$ and $F_c$ (move j from $F_c$ to $E_c$)*

9. *Move to next partition: $\alpha_c \leftarrow \alpha_{next}$*

10. *For each sub-partition $\alpha_c^i\epsilon\alpha_c$ such that i) $|\alpha_c^i| > 1$; ii) $\alpha_Y|_{\alpha_c^i}$ is not yet classified; and iii) $F_c$ is not empty*

    start the **Search** procedure (for sub-partition $\alpha_c^i$, $F_c$, $E_c$)

*If $\alpha_Y|_{\alpha_c}$ is classified, return.*

*If $|\alpha_c| = 1$ classify according to the instance's class value, return.*

*If $F_c = \{\emptyset\}$ classify according to the most common value of the class attribute, return.*

# The Proposed DID Algorithm

Function Search: Given a sub-partition $\alpha^i_c$ and the attribute partitions remaining $F_c \subseteq F$

1. Init current partition $\alpha_c$ by $\alpha^i_c$; init $E_c \leftarrow E$; $F_c \leftarrow F$

2. Normalize probabilities of the elements of $\alpha^i_c$.

3. Create local class partition $\alpha_Y|_{\alpha_c}$.

4. Generate neighborhood partitions:

   $N(\alpha_c) = \{\alpha_j|_{\alpha_c}\}, j \in F_c, F_c = \{j : A_j \text{ not selected by the algorithm yet}\}$

5. Normalize probabilities of neighbors and of the class partition;

resolved arbitrary):

$$\alpha_{next} \leftarrow \arg\min_{\alpha_j \in N(\alpha_c)}\{w_1 d_1(\alpha_c, \alpha_j|_{\alpha_c}) + w_2 d_2(\alpha_j|_{\alpha_c}, \alpha_Y|_{\alpha_c})\}$$

8. Update $E_c$ and $F_c$ (move $j$ from $F_c$ to $E_c$)

9. Move to next partition: $\alpha_c \leftarrow \alpha_{next}$

10. For each sub-partition $\alpha^i_c \in \alpha_c$ such that i) $|\alpha^i_c| > 1$; ii) $\alpha_Y|_{\alpha^i_c}$ is not yet

    classified; and iii) $F_c$ is not empty

       start the **Search** procedure (for sub-partition $\alpha^i_c$, $F_c$, $E_c$)

If $\alpha_Y|_{\alpha_c}$ is classified, return.

If $|\alpha_c| = 1$ classify according to the instance's class value, return.

If $F_c = \{\emptyset\}$ classify according to the most common value of the class attribute,

   return.

# The Proposed DID Algorithm

**Function Search.** *Given set $\alpha_c^i$ and the attributes partitions $\alpha_1, \alpha_2, \ldots, \alpha_n$; $E$; $F$;*

1. *Init current partition $\alpha_c$ by $\alpha_c^i$; init $E_c \leftarrow E$; $F_c \leftarrow F$*

2. *Normalize probabilities of the elements of $\alpha_c^i$.*

3. *Create local class partition $\alpha_Y|_{\alpha_c}$.*

4. *Generate neighborhood partitions:*

   $N(\alpha_c) = \{\alpha_c|\ldots\}, j \in F_c, F_c = \{j : A_j$ not selected by the algorithm yet$\}$

6. *Obtain distance measures by $d_1(\alpha_c, \alpha_j|_{\alpha_c})$ and $d_2(\alpha_j|_{\alpha_c}, \alpha_Y|_{\alpha_c}), j \in F_c$*

7. *Choose next partition:* The next partition is selected as follows (ties are resolved arbitrary):

$$\alpha_{next} \leftarrow \arg\min_{\alpha_j \in N(\alpha_c)} \{w_1 d_1(\alpha_c, \alpha_j|_{\alpha_c}) + w_2 d_2(\alpha_j|_{\alpha_c}, \alpha_Y|_{\alpha_c})\}$$

8. *Update $E_c$ and $F_c$ (move $j$ from $F_c$ to $E_c$)*

9. *Move to next partition: $\alpha_c \leftarrow \alpha_{next}$*

10. *For each sub-partition $\alpha_c^i \in \alpha_c$ such that i) $|\alpha_c^i| > 1$; ii) $\alpha_Y|_{\alpha_c^i}$ is not yet classified; and iii) $F_c$ is not empty*

    *start the* **Search** *procedure (for sub-partition $\alpha_c^i$, $F_c$, $E_c$)*

*If $\alpha_Y|_{\alpha_c}$ is classified, return.*

*If $|\alpha_c| = 1$ classify according to the instance's class value, return.*

*If $F_c = \{\emptyset\}$ classify according to the most common value of the class attribute, return.*

# The Proposed DID Algorithm

**Function Search.** *Given set $\alpha_c^i$ and the attributes partitions $\alpha_1, \alpha_2, \dots, \alpha_n$; $E$; $F$:*

1. *Init current partition $\alpha_c$ by $\alpha_c^i$; init $E_c \leftarrow E$; $F_c \leftarrow F$*

2. *Normalize probabilities of the elements of $\alpha_c^i$.*

3. *Create local class partition $\alpha_Y|_{\alpha_c}$.*

4. *Generate neighborhood partitions:*

   $N(\alpha_c) = \{\alpha_j|_{\alpha_c}\}, j \epsilon F_c, F_c = \{j : A_j$ not selected by the algorithm yet$\}$

5. *Normalize probabilities of neighbors and of the class partition;*

6. *Obtain distance measures by $d_1(\alpha_c, \alpha_j|_{\alpha_c})$ and $d_2(\alpha_j|_{\alpha_c}, \alpha_Y|_{\alpha_c}), j \epsilon F_c$*

7. *Choose next partition: The next partition is selected as follows (ties are resolved arbitrary):*

8. *Update $E_c$ and $F_c$ (move $j$ from $F_c$ to $E_c$)*

9. *Move to next partition: $\alpha_c \leftarrow \alpha_{next}$*

10. *For each sub-partition $\alpha_c^i \epsilon \alpha_c$ such that i) $|\alpha_c^i| > 1$; ii) $\alpha_Y|_{\alpha_c^i}$ is not yet classified; and iii) $F_c$ is not empty*

    *start the* **Search** *procedure (for sub-partition $\alpha_c^i$, $F_c$, $E_c$)*

*If $\alpha_Y|_{\alpha_c}$ is classified, return.*

*If $|\alpha_c| = 1$ classify according to the instance's class value, return.*

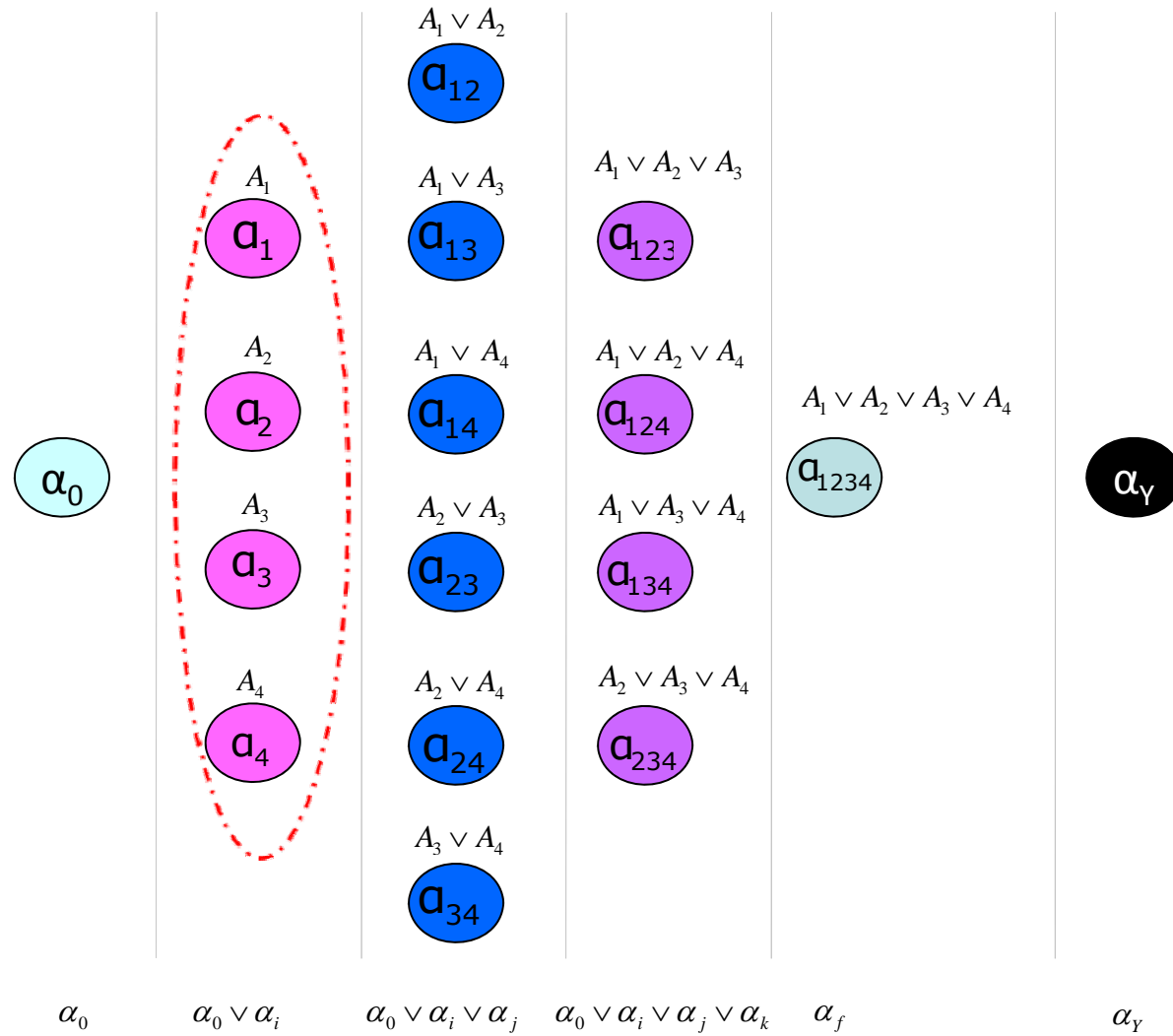*If $F_c = \{\emptyset\}$ classify according to the most common value of the class attribute, return.*

# The Proposed DID Algorithm

**Function Search.** *Given set $\alpha_c^i$ and the attributes partitions $\alpha_1, \alpha_2, \ldots, \alpha_n$; $E; F$;*

1. *Init current partition $\alpha_c$ by $\alpha_c^i$; init $E_c \leftarrow E$; $F_c \leftarrow F$*

2. *Normalize probabilities of the elements of $\alpha_c^i$.*

3. *Create local class partition $\alpha_Y|_{\alpha_c}$.*

4. *Generate neighborhood partitions:*

   $N(\alpha_c) = \{\alpha_j|_{\alpha_c}\}, j \epsilon F_c, F_c = \{j : A_j$ not selected by the algorithm yet$\}$

5. *Normalize probabilities of neighbors and of the class partition;*

6. *Obtain distance measures by $d_1\left(\alpha_c, \alpha_j|_{\alpha_c}\right)$ and $d_2\left(\alpha_j|_{\alpha_c}, \alpha_Y|_{\alpha_c}\right), j \epsilon F_c$*

7. *Choose next partition: The next partition is selected as follows (ties are resolved arbitrary):*

   $\alpha_{next} \leftarrow arg\,min_{\alpha_j \in N(\alpha_c)}\{w_1 d_1\left(\alpha_c, \alpha_j|_{\alpha_c}\right) + w_2 d_2\left(\alpha_j|_{\alpha_c}, \alpha_Y|_{\alpha_c}\right)\}$

10. *For each sub-partition $\alpha_c^i \epsilon \alpha_c$ such that i) $\left|\alpha_c^i\right| > 1$; ii) $\alpha_Y|_{\alpha_c^i}$ is not yet classified; and iii) $F_c$ is not empty*

    *start the **Search** procedure (for sub-partition $\alpha_c^i$, $F_c$, $E_c$)*

*If $\alpha_Y|_{\alpha_c}$ is classified, return.*

*If $\left|\alpha_c\right| = 1$ classify according to the instance's class value, return.*

*If $F_c = \{\emptyset\}$ classify according to the most common value of the class attribute, return.*

# Outline

1. Introduction & Motivation

2. Proposed Partitions Approach

3. **Examples**

4. Results

5. Mid-level solutions
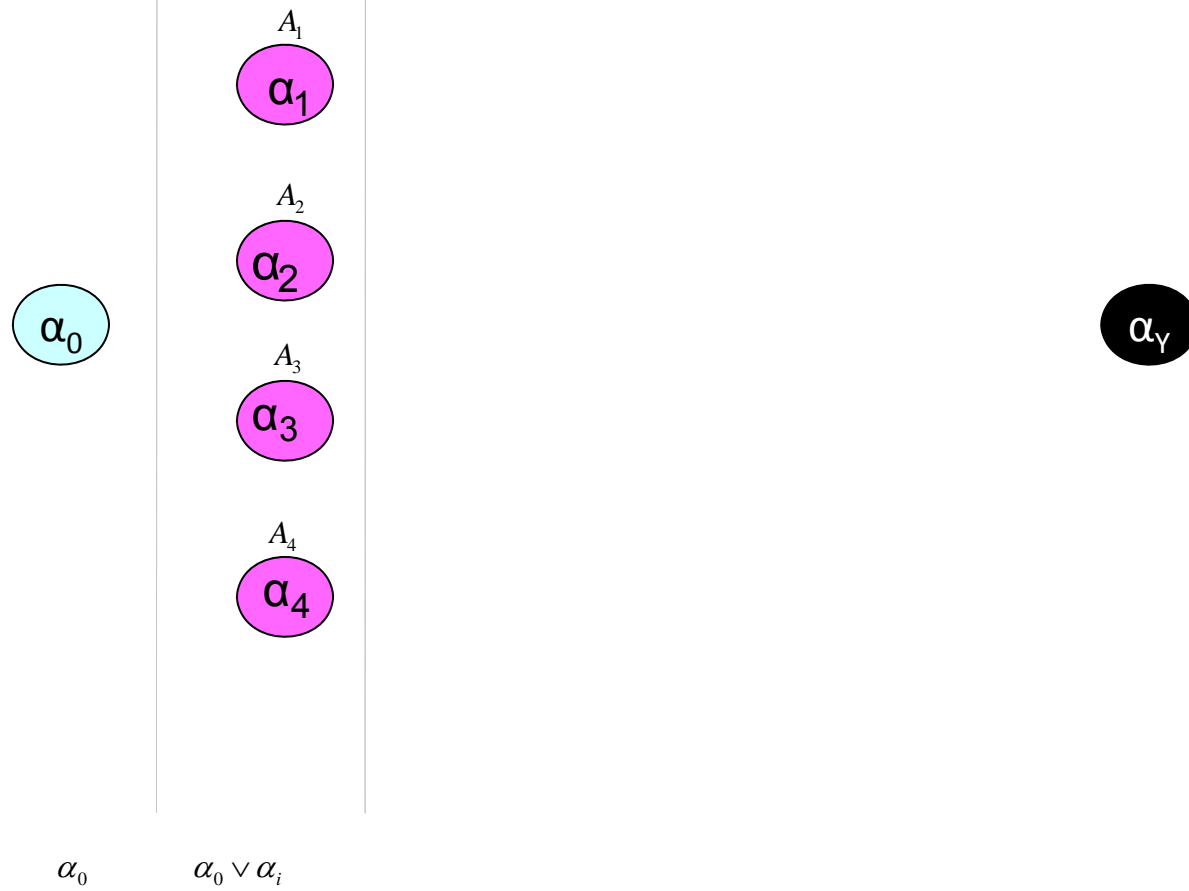
6. Summary & Contribution

# Example

**Training Data Set**

|    | Y | A1 | A2 | A3 | A4 |
|----|---|----|----|----|----|
| 1  | 1 | 4  | 1  | 1  | 1  |
| 2  | 2 | 1  | 1  | 2  | 2  |
| 3  | 2 | 1  | 1  | 2  | 1  |
| 4  | 2 | 1  | 1  | 2  | 1  |
| 5  | 3 | 2  | 3  | 2  | 1  |
| 6  | 3 | 2  | 2  | 1  | 1  |
| 7  | 4 | 2  | 3  | 1  | 1  |
| 8  | 4 | 3  | 3  | 1  | 1  |
| 9  | 5 | 3  | 3  | 3  | 1  |
| 10 | 4 | 1  | 1  | 3  | 1  |
| 11 | 4 | 1  | 2  | 4  | 2  |
| 12 | 4 | 2  | 2  | 4  | 2  |

ID3 example

$A_1 \vee A_2$

$\alpha_{12}$

$A_1$

$\alpha_1$

$A_1 \vee A_3$

$\alpha_{13}$

$A_1 \vee A_2 \vee A_3$

$\alpha_{123}$

$A_2$

$\alpha_2$

$A_1 \vee A_4$

$\alpha_{14}$

$A_1 \vee A_2 \vee A_4$

$\alpha_{124}$

$A_1 \vee A_2 \vee A_3 \vee A_4$

$\alpha_0$

$\alpha_{1234}$

$\alpha_Y$

$A_3$

$\alpha_3$

$A_2 \vee A_3$

$\alpha_{23}$

$A_1 \vee A_3 \vee A_4$

$\alpha_{134}$

$A_4$

$\alpha_4$

$A_2 \vee A_4$

$\alpha_{24}$

$A_2 \vee A_3 \vee A_4$

$\alpha_{234}$

$A_3 \vee A_4$

$\alpha_{34}$

$\alpha_0$     $\alpha_0 \vee \alpha_i$     $\alpha_0 \vee \alpha_i \vee \alpha_j$     $\alpha_0 \vee \alpha_i \vee \alpha_j \vee \alpha_k$     $\alpha_f$     $\alpha_Y$

## ID3 example

$\alpha_1$    $A_1$

$\alpha_2$    $A_2$

$\alpha_0$

$\alpha_3$    $A_3$

$\alpha_Y$

$\alpha_4$    $A_4$

$\alpha_0$      $\alpha_0 \vee \alpha_i$

# ID3 example

## Classification Tree

$A_1$

$A_1=1$   $A_1=2$   $A_1=3$   $A_1=4$

## Partitions Graph

$\alpha_1$

0.904

$\alpha_2$   1.3

$\alpha_0$   $\alpha_Y$

0.937

$\alpha_3$

1.87

$\alpha_4$

Results for :

$w_1 = 0; w_2 = +1;$

distance2 = Conditional Entropy

# ID3 example

## Classification Tree



$A_1$

$A_1=1$

$A_3$

## Partitions Graph



$\alpha_{12}$

$\alpha_1$

$\alpha_{13}$

$\alpha_{14}$

0.65

0

0.95

$\alpha_Y$

**Restricted to $A_1=1$ subset**

$\alpha_{23}$

$\alpha_{24}$

$\alpha_{34}$

Results for :

$w_1 = 0; w_2 = +1;$

distance2 = Conditional Entropy

# ID3 example

**Classification Tree**

**Partitions Graph**

$A_1$

$A_1=1$  $A_1=2$

$A_3$  $A_3$

$\alpha_1$

$\alpha_{12}$

$\alpha_{13}$

$\alpha_{14}$

$\alpha_{23}$

$\alpha_{24}$

$\alpha_{34}$

$\alpha_Y$

1

0.5

0.69

**Restricted to $A_1=2$ subset**

Results for :

$w_1 = 0$; $w_2 = +1$ ;

distance2 = Conditional Entropy

# ID3 example

## Classification Tree

## Partitions Graph

Restricted to $A_1=3$ subset

Results for :

$w_1 = 0$; $w_2 = +1$;

distance2 = Conditional Entropy

# ID3 example

## Classification Tree

## Partitions Graph

# ID3 example

## Classification Tree

Average depth = 2.1

No. of decision = 5

No. of leaves = 10

Max steps = 3

## Partitions Graph

Results for :
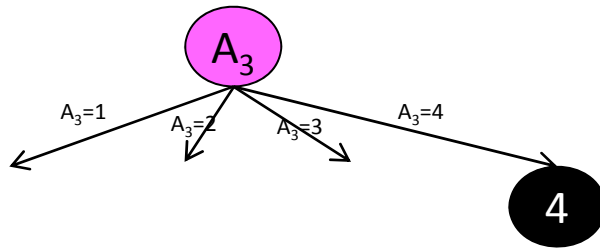$w_1 = 0$; $w_2 = +1$;
distance2 = Conditional Entropy

# The DID approach

$A_1 \vee A_2$

$A_{12}$

$A_1 \vee A_3$

$A_1$

$\alpha_1$

$A_{13}$

$A_1 \vee A_2 \vee A_3$

$\alpha_{123}$

$A_2$

$\alpha_2$

$A_1 \vee A_4$

$A_{14}$

$A_1 \vee A_2 \vee A_4$

$\alpha_{124}$

$A_1 \vee A_2 \vee A_3 \vee A_4$

$\alpha_0$

$\alpha_{1234}$

$\alpha_Y$

$A_3$

$\alpha_3$

$A_2 \vee A_3$

$A_{23}$

$A_1 \vee A_3 \vee A_4$

$\alpha_{134}$

$A_4$

$\alpha_4$

$A_2 \vee A_4$

$A_{24}$

$A_2 \vee A_3 \vee A_4$

$\alpha_{234}$

$A_3 \vee A_4$

$A_{34}$

$\alpha_0$  $\quad\quad$  $\alpha_0 \vee \alpha_i$  $\quad\quad$  $\alpha_0 \vee \alpha_i \vee \alpha_j$  $\quad\quad$  $\alpha_0 \vee \alpha_i \vee \alpha_j \vee \alpha_k$  $\quad\quad$  $\alpha_f$  $\quad\quad$  $\alpha_Y$

# The DID approach

$A_1$

$\alpha_1$

$A_2$

$\alpha_2$

$\alpha_0$

$A_3$

$\alpha_3$

$\alpha_Y$

$A_4$

$\alpha_4$

$\alpha_0$ $\qquad$ $\alpha_0 \vee \alpha_i$

# The DID approach

**Classification Tree**



**Partitions Graph**



Results for :

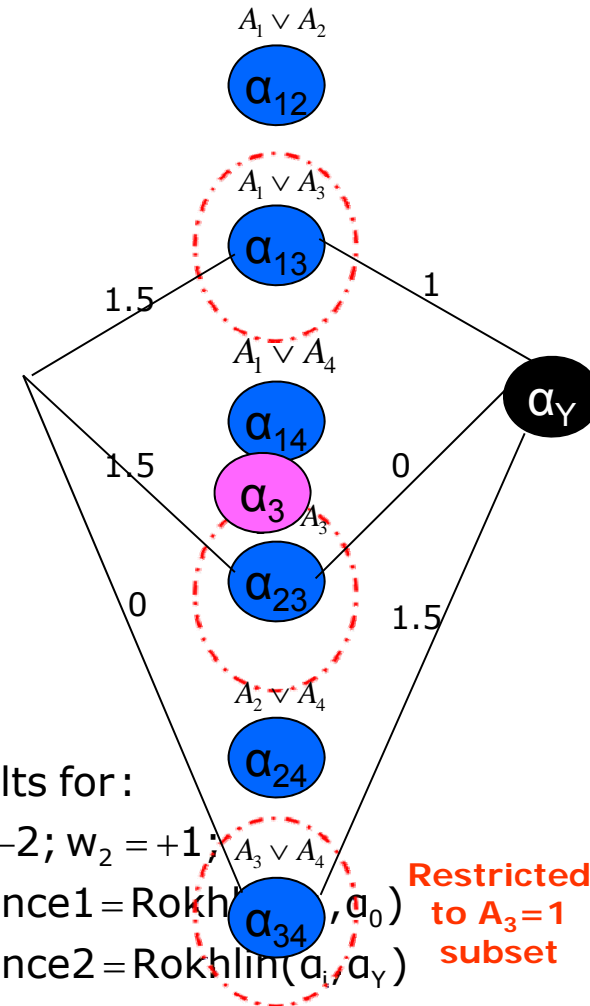$w_1 = -2; w_2 = +1;$

$distance1 = Rokhlin(\alpha_i, \alpha_0)$

$distance2 = Rokhlin(\alpha_i, \alpha_Y)$

# The DID approach

**Classification Tree**

**Partitions Graph**



$A_1 \vee A_2$

$\alpha_{12}$

$A_3$

$A_3=1$  $A_3=2$  $A_3=3$  $A_3=4$

$A_1 \vee A_3$

$\alpha_{13}$

$A_2$

4

1.5                1

1

2

3

$A_1 \vee A_4$

1

3

4

$\alpha_{14}$

$\alpha_Y$

1.5              $\alpha_3$          0

$A_3$

$\alpha_{23}$

0                              1.5

$A_2 \vee A_4$

$\alpha_{24}$

Results for :
$w_1 = -2; w_2 = +1;$   $A_3 \vee A_4$   **Restricted**
distance1 = Rokhlin$(\alpha_i, \alpha_0)$   $\alpha_{34}$   **to $A_3=1$**
distance2 = Rokhlin$(\alpha_i, \alpha_Y)$   **subset**

# Some Results

# Summarizing Comparison between ID3, C4.5 and DID decision trees

| Dataset | Size | | ID3 | | C4.5 | | DID | |
|---|---|---|---|---|---|---|---|---|
| | #instances | #Attributes | Average Depth | Accuracy | Average Depth | Accuracy | Average Depth | Accuracy |
| Monk's-1 | 124 | 6 | 3.21 | 82% | 3.32 | 82% | 2.66 | 96.% |
| Monk's-2 | 169 | 6 | 4.34 | 70.4% | 4.6 | 75% | 4.2 | 66% |
| Monk's full Random set | 216 | 6 | 1.93 | 100% | 2.04 | 100% | 1.8 | 100% |
| Connect4 | 67,557 | 42 | 5.85 | 73.8% | 10.16 | 79.4% | 5.64 | 75% |
| SPECT Heart | 80 | 22 | 9.6 | 75.1% | 10.2 | 80.3% | 9.3 | 76% |
| Voting | 435 | 16 | 1.8 | 96% | 2.2 | 96.6% | 2.1 | 96% |
| Balance Scale | 625 | 4 | 3.4 | 76.3% | 3.4 | 78.6% | 3.3 | 76.6% |
| Cars | 1728 | 6 | 2.82 | 77.1% | 2.83 | 77% | 2.77 | 78.5% |
| Tic-Tac-Toe | 958 | 9 | 4.62 | 80.6% | 4.62 | 80.4% | 4.6 | 76.2% |
| Soy Beans | 47 | 35 | 1.35 | 100% | 2.37 | 97% | 1.32 | 97% |
| Lymphography | 148 | 18 | 2.71 | 75.1% | 6.51 | 77.3% | 2.6 | 72.6% |

| Case | #features | SVM accuracy% | J48 % | DID accuracy % |
|---|---|---|---|---|
| australian | 14 | 55.5 | 86.2 | **86.9** |
| breast | 9 | **96.5** | 93.6 | 93.5 |
| diabetes | 8 | 65.1 | **74.2** | 72.6 |
| glass | 8 | **69.16** | 50.1 | 51.8 |
| glass2 | 8 | 76.68 | 75.3 | **82.1** |
| heart | 13 | 55.93 | 79.4 | **79.5** |
| iris | 4 | **96.67** | 94.4 | 95.6 |
| pima | 36 | 65.1 | **73.1** | 72.2 |
| segment | 18 | 63.9 | **94.1** | 93.6 |
| Shuttle-small | 9 | **89.41** | 62 | 61.9 |
| vehicle | 18 | 30.5 | **69.7** | 65.2 |
| waveform-21 | 21 | **86.1** | 76.3 | 73.5 |
| cleve | 13 | 54.73 | **78.9** | 78 |
| crx | 15 | 65.67 | 87.5 | **87.6** |
| german | 20 | **70** | 65.2 | 66.6 |
| hepatitis | 19 | **83.55** | 57.4 | 64.2 |
| chess | 36 | 93.83 | 99.3 | **99.8** |
| corral | 6 | 96.89 | 98.1 | **98.3** |
| flare | 9 | **82.37** | 61.2 | 68.9 |
| mofn-3-7-10 | 10 | **100** | **100** | **100** |
| soybean-large | 35 | 87.19 | **95.8** | 94.2 |
| vote | 15 | 95.35 | 94.7 | **95.4** |

# Accuracy of C4.5 (J48) and DID as a function of the number of removed features for different cases taken from the UCI Repository

# Outline

1. Introduction & Motivation

2. Our Partitions Approach

3. Example

4. Results

5. Mid-level solutions

6. <span style="color:orange">Summary & Contribution</span>

## Summary & Contribution

**Modeling the tree construction problem as a shortest path problem over a graph of partitions as nodes.**

- ❑ A <u>unified framework</u> for existing DT algorithms
- ❑ Further Generalization <u>via different metrics</u>, e.g, Rokhlin, Entropy, etc. supported by IT
- ❑ Orthogonally vs. Information Gain
- ❑ Big Data fit: Shorter trees with smaller decisions for online scoring and recommendation

# Thank you !
## Questions?