

# When Global Optimization Meets Machine Learning: Challenges and Opportunities

Yankai Cao

University of British Columbia



THE UNIVERSITY  
OF BRITISH COLUMBIA

# Machine Learning vs Human Learning

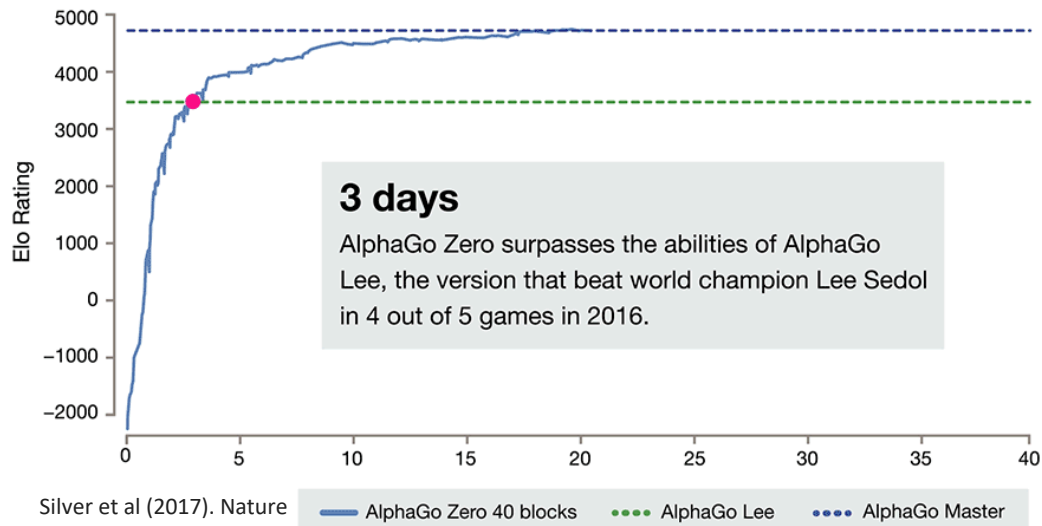


**Why are humans more data-efficient?**

4.9 million games vs 3 thousand

**Hypothesis:**

- Transfer learning
- Reasoning + Hybrid Modelling
- Higher solution accuracy?  
→ Can global optimization improve the performance of ML?



# Prevailing Mindsets in Machine Learning Community

The most obvious drawback of the learning procedure is that the error-surface may contain **local minima** so that gradient descent is not guaranteed to find a global minimum. However, experience with many tasks shows that the network very rarely gets stuck in poor local minima that are significantly worse than the global minimum.

— David Rumelhart, **Geoffery Hinton**, Ronald Williams  
*Learning Representations by Back-propagating Errors*, Nature (1986)

In the late 1990s, neural nets and backpropagation were largely **forsaken** by the machine-learning community and ignored by the computer-vision and speech-recognition communities. It was widely thought that learning useful, multistage, feature extractors with little prior knowledge was infeasible. In particular, it was commonly thought that simple gradient descent would get trapped in **poor local minima** ...

— Yann LeCun, Yoshua Bengio, **Geoffery Hinton**  
*Deep Learning*, Nature (2015)

# Prevailing Mindsets in Machine Learning Community – Part 2

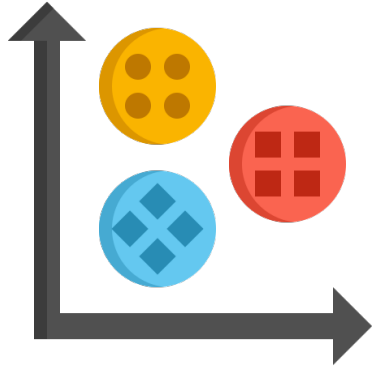
Belief 1 (1990s - now): Solving machine learning problems to global optimum is impossible using state-of-the-art solvers, especially for large datasets.

- It is widely believed that deterministic global algorithms are too slow (Hamm et al., 2007).
- In 1990s and 2000s, stochastic global optimization methods were very popular, including Evolving Neural Networks (Yao et al., 1999, Stanley et al., 2002), Genetic Algorithms (Miller et al., 1989, Leung et al., 2003), and Simulated Annealing (Goffe et al., 1994).

Belief 2 (2010s - now): Solving ML problem to local optimum is good enough.

- With strong simplification assumptions on the distribution of data and network weight parameters, as the number of hidden units increases, all local minima become increasingly close to being global minima (Choromanska et al. 2015).
- If a neural net is strongly over-parameterized so that I can memorize any dataset, then all critical points become global minima (Livni et al. 2014, Nguyen & Hein 2017).
- For neural nets with one hidden layer and a convolutional structure with no overlap, and a ReLU activation function, if inputs are independent gaussian variables, gradient descent converges to the global optimum (Bruzkus & Globerson, 2017)

# We Focus on the Global Optimization of

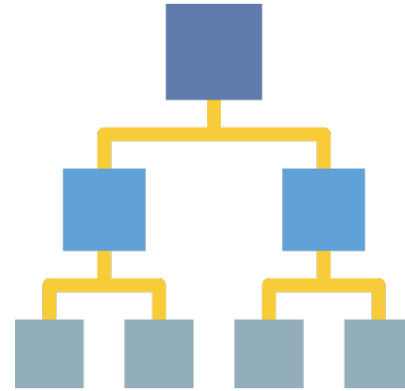


Centroid-based Clustering

[Hua et al. ICML 2021<sup>spotlight</sup>]

[Shi et al. ICML 2022]

[Ren et al. NeurIPS 2022]



Decision Tree

[Hua et al. NeurIPS 2022]

2022 INFORMS Data Mining Best Paper

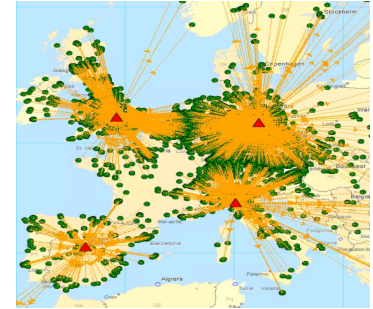
# Applications with High Stake



Bail Decision



Medical Diagnosis



Facility Location Problem

Interpretable

Intuitive

?

Heuristic Algorithm -> Suboptimal Solutions

# Our Hypothesis

## Deterministic Global Optimality

- Solving ML problem to global optimum is possible for **large datasets** (1 million samples).



## Enhanced Solution

- The global optimal solution renders a **much better model** than the local solution ( $\geq 1\%$  improvement).

# Outline

- ML as Stochastic Mixed Integer Problems
  - Reduced Space Branch and Bound Algorithm
- Centroid-based Clustering
- Optimal Decision Tree





# Machine Learning Problems are MI(N)LP Problems

Given:

- $\{x_1, \dots, x_s, \dots, x_n\} \in \mathbb{R}^P$ ,
- $\{y_1, \dots, y_s, \dots, y_n\} \in \{0,1\}^K$ ,
- $\mathcal{S} = \{1, \dots, s, \dots, n\}$
- Model representation  $h$

Output:

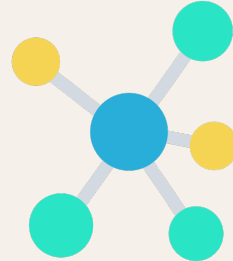
- Optimal parameters  $m \in M$

Branch and Bound (BB)  
Framework



Supervised Learning

$$\begin{aligned} \mathcal{L}(M) &:= \min_{m \in M, \hat{y}_s} \sum_{s \in \mathcal{S}} E(\hat{y}_s, y_s) + \gamma R(m) \\ \text{s. t.} & \quad \hat{y}_s = h(m, x_s) \\ & \quad s \in \mathcal{S} \end{aligned}$$



Unsupervised Learning

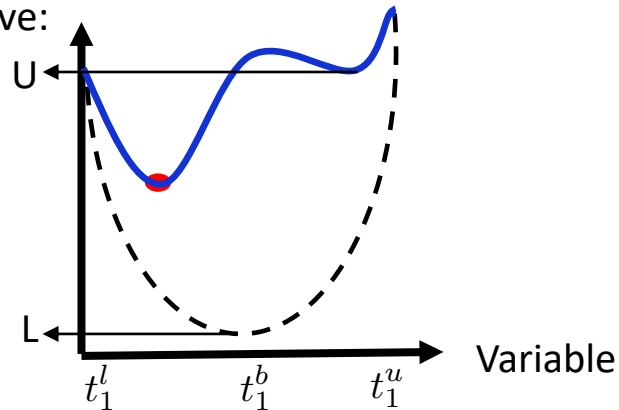
$$\begin{aligned} \mathcal{L}(M) &:= \min_{m \in M, \hat{x}_s} \sum_{s \in \mathcal{S}} E(\hat{x}_s, x_s) + \gamma R(m) \\ \text{s. t.} & \quad \hat{x}_s = h(m, x_s) \\ & \quad s \in \mathcal{S} \end{aligned}$$

# Branch and Bound Framework

$$\mathcal{L}(M) := \min_{m \in M, \hat{y}_s} \sum_{s \in \mathcal{S}} E(\hat{y}_s, y_s) + \gamma R(m)$$

*s.t.*      $\hat{y}_s = h(m, x_s)$   
               $s \in \mathcal{S}$

Objective:



Root  
node

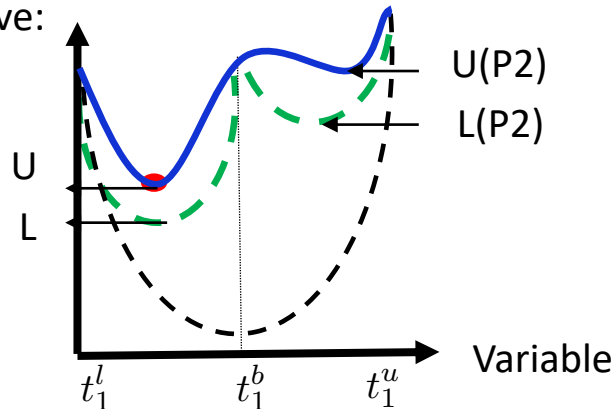


# Branch and Bound Framework

$$\mathcal{L}(M) := \min_{m \in M, \hat{y}_s} \sum_{s \in \mathcal{S}} E(\hat{y}_s, y_s) + \gamma R(m)$$

s.t.  $\hat{y}_s = h(m, x_s)$   
 $s \in \mathcal{S}$

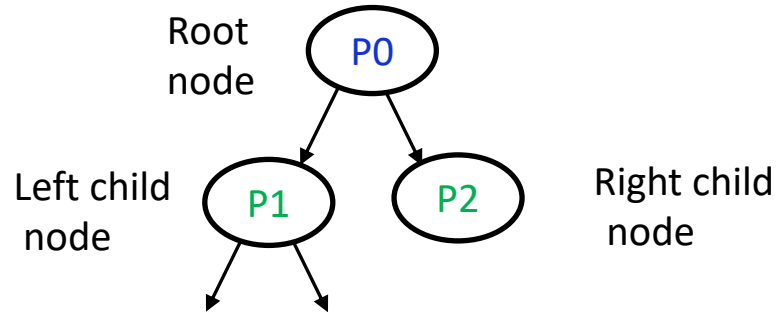
Objective:



**Limitation:** Need to branch on the **full variable space**

**For Decision Tree:** # Integer  $2^D(n + P + K + 1)$   
e.g. depth  $D = 2$  and  $n = 1$  million samples  
-> **4 million integer variables**

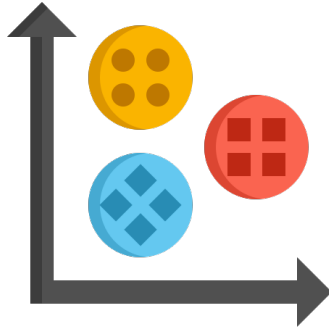
**Key:** Need to exploit the structure



# Machine Learning Problems are Two-stage Stochastic Programming Problems

First Stage  
Variables

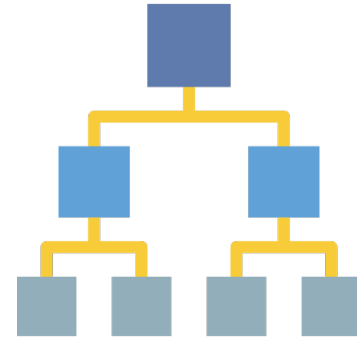
Model  
Parameter



Centroid-based Clustering

Second Stage  
Variables

Sample  
Evaluation



Decision Tree

# Machine Learning Problems are Two-stage Stochastic Programming Problems

First Stage  
Variables

Model  
Parameter

Second Stage  
Variables

Sample  
Evaluation

$$\mathcal{L}(M) := \min_{m \in M, \hat{x}_s} \sum_{s \in \mathcal{S}} E(\hat{x}_s, x_s) + \gamma R(m)$$

s. t.  $\hat{x}_s = h(m, x_s)$   
 $s \in \mathcal{S}$



Where:

$$\mathcal{L}(M) := \min_{m \in M} \sum_{s \in \mathcal{S}} Q_s(m)$$
$$Q_s(m) := \min_{\hat{y}_s} E(\hat{y}_s, x_s) + \frac{\gamma}{|\mathcal{S}|} R(m)$$

s. t.  $\hat{y}_s = h(m, x_s)$

Various approaches proposed in the stochastic programming community can be utilized to exploit the problem structure (Karupiah & Grossmann 2008, Li et al. 2011, Cao & Zavala 2019, Li & Grossman 2019, Ogbe & Li 2019, Kannan 2019).

# Reduced Space Branch and Bound

**Machine Learning Problems:**

$$\mathcal{L}(M) := \min_{m \in M} \sum_{s \in \mathcal{S}} Q_s(m)$$

**Lower Bounding Problem:**  
(wait and see problem)

$$\beta(M) := \min_{m_s \in M} \sum_{s \in \mathcal{S}} Q_s(m_s)$$

**Upper Bounding Problem:**  
(fix  $\hat{m} \in M$ )

$$\alpha(M) := \sum_{s \in \mathcal{S}} Q_s(\hat{m})$$

## Assumption

Around solution  $m^*$ ,  $Q_s(m)$  is continuous with respect to  $m_c$  (i.e., continuous first-stage).

## Theorem [Cao & Zavala, 2019]

With the above lower and upper bounds, the BB procedure can *converge* by branching **only on first stage variables** (i.e., model parameters)

**For Decision Tree:** e.g. depth  $D = 2$  and  $n = 1$  million samples, with 3 features and 2 classes

branching variables: 4 million  $\rightarrow$  23

Independent to the number of samples

# Reduced Space Branch and Bound

Lower Bound Strategy

Upper Bound Strategy

Bound Tightening

Centroid-based  
Clustering

Decision Tree

# Clustering Problems

Given: Dataset  $X = \{x_1, \dots, x_S\} \in \mathbb{R}^{A \times S}$

Desired cluster number,  $K$

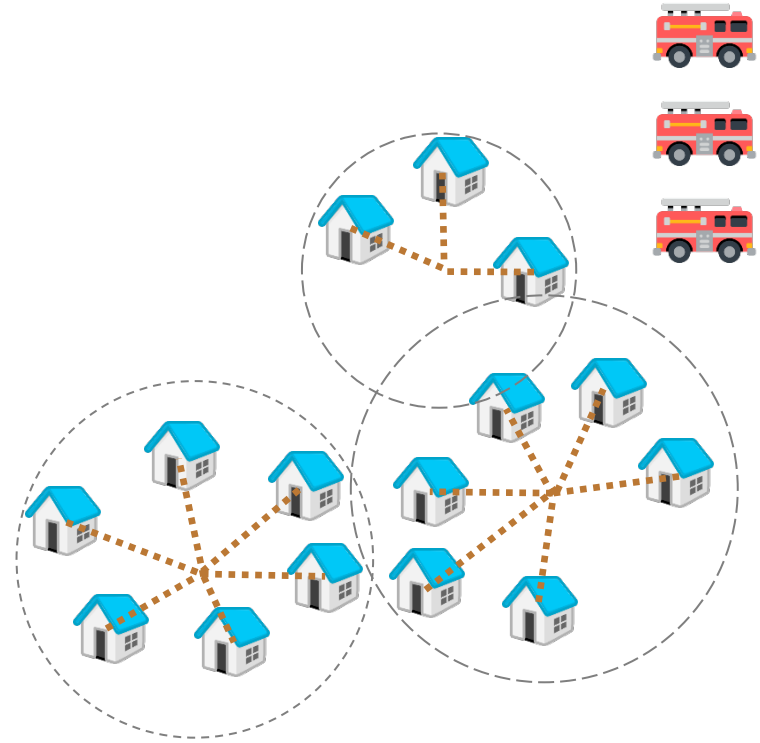
Output:  $m_k$  the center of cluster  $k$

Object:

$$k\text{-means} : \min_m \sum_{s \in S} \min_{k \in \mathcal{K}} \|x_s - m_k\|_2^2$$

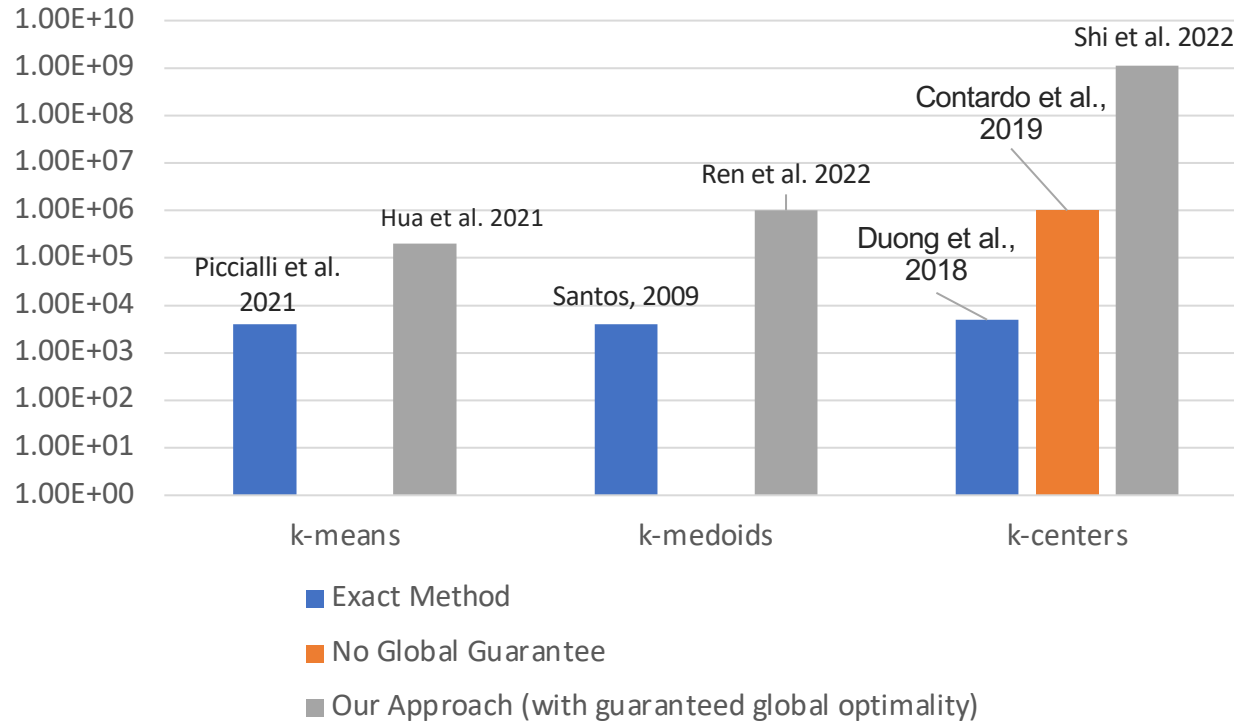
$$k\text{-medoids} : \min_{m \in X} \sum_{s \in S} \min_{k \in \mathcal{K}} \|x_s - m_k\|_2^2$$

$$k\text{-centers} : \min_{m \in X} \max_{s \in S} \min_{k \in \mathcal{K}} \|x_s - m_k\|_2^2$$





# Recent Progress on the Scalability of Data Size



# MINLP Formulation

$$k\text{-means} : \min_m \sum_{s \in \mathcal{S}} \min_{k \in \mathcal{K}} \|x_s - m_k\|_2^2$$

$$b_{s,k} = \begin{cases} 1 & \text{if } x_s \text{ is in cluster } k. \\ 0 & \text{otherwise.} \end{cases}$$

$d_{s,k}$  is the distance between  $x_s$  and the cluster center  $m_k$

$d_{s,*}$  is  $\min_k d_{s,k}$

First Stage Variables

$$m: K \times P$$

(Center of clusters, continuous)

$$\min_{m,d,b} \sum_{s \in \mathcal{S}} d_{s,*}$$

$$s. t. -N(1 - b_{s,k}) \leq d_{s,*} - d_{s,k} \leq N(1 - b_{s,k})$$

$$d_{s,k} \geq \|x_s - m_k\|^2$$

$$\sum_{k \in \mathcal{K}} b_{s,k} = 1$$

$$b_{s,k} \in \{0,1\}$$

$$s \in \mathcal{S}, k \in \mathcal{K}$$

Second Stage Variables

$$b_{s,k}, d_{s,k}, d_{s,*}$$

(Cluster assignment, mixed integer)

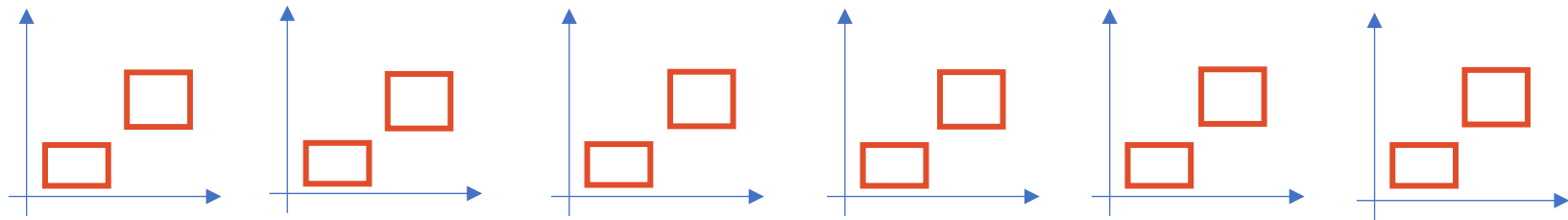
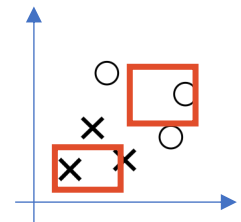
# Basic Lower Bound

$$\min_{m \in M} \sum_{s \in \mathcal{S}} Q_s(m) \quad \longrightarrow \quad \beta^B(M) := \min_{m_s \in M} \sum_{s \in \mathcal{S}} Q_s(m_s) \quad \longrightarrow \quad \beta_S(M) := \min_{m_s \in M} Q_s(m_s)$$

s. t.  $m_s = m_{s+1}$

non-anticipativity constraint

Decompose into  $|\mathcal{S}| = n$  subproblems



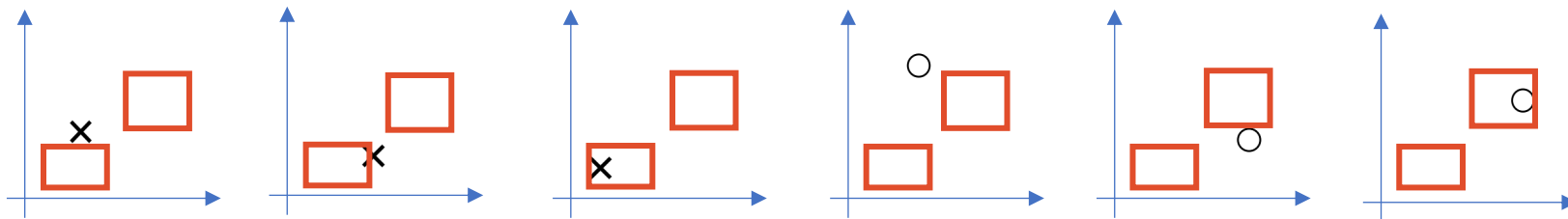
# Basic Lower Bound

$$\beta_S(M) := \min_{m_S \in M} Q_S(m_S) \quad \longrightarrow \quad \beta_S(M) := \min_k \beta_{S,k}(M_k) = \min_k \min_{m_{S,k} \in M_k} \|x_S - m_{S,k}\|_2^2$$

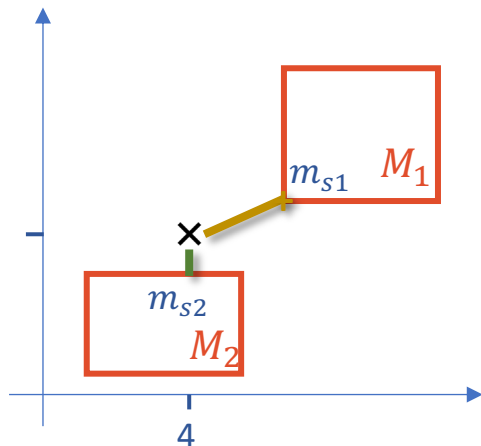
$M_k := \{m_{S,k} \mid \underline{m}_k \leq m_{S,k} \leq \bar{m}_k\}$

Decompose into  $|\mathcal{S}|$   
=  $n$  subproblems

$\beta_{S,k}$  has **closed form** solution:  
 $m_{S,k,i} = \text{med}\{\underline{m}_{k,i}, x_{S,i}, \bar{m}_{k,i}\}, i \in \{1, \dots, P\}$



# Basic Lower Bound



$$M_k := \{m_{s,k} \mid \underline{m}_k \leq m_{s,k} \leq \bar{m}_k\}$$

$$\beta_s(M) := \min_k \beta_{s,k}(M_k) = \min_k \min_{m_{s,k} \in M_k} \|x_s - m_{s,k}\|_2^2$$

$\beta_{s,k}$  has **closed form** solution:

$$m_{s,k,i} = \text{med}\{\underline{m}_{k,i}, x_{s,i}, \bar{m}_{k,i}\}, i \in \{1, \dots, P\}$$

Assuming 2 cluster, 2 features ( $K=2, P=2$ ):

Subproblem  $s$ :  $x_s = (4, 5)$        $\beta_s(M) = 1$

- Cluster 1 center:  $m_{s,1,1} = \text{med}\{4, 6, 10\} = 6$

$$m_{s,1,2} = \text{med}\{5, 6, 10\} = 6$$

- Object value:  $\beta_{s,1}(M_1) = \sqrt{(6-4)^2 + (6-5)^2} = \sqrt{5}$

Calculation can be  
easily parallelized

No need to solve any optimization  
problem

# Lagrangian Relaxation

Dualized the non-anticipativity constraints and added to the objective functions with Lagrange multipliers  $\lambda$

$\lambda$  can be optimized via sub-gradient method [Held and Karp, 1971]

**Proposition:**  
 $\beta^B(M) = \beta^{LD}(M, 0) \leq \beta^{LD}(M) \leq \mathcal{L}(M)$

$$\beta^{LD}(M, \lambda) := \min_{m_s \in M} \sum_{s \in \mathcal{S}} Q_s(m_s) + \sum_{s=1}^{|\mathcal{S}|-1} \lambda_s^T (m_s - m_{s+1})$$



$$\beta_s^{LD}(M, \lambda) := \min_{m_s \in M} Q_s(m_s) + (\lambda_s - \lambda_{s+1})^T m_s$$



$$\beta^{LD}(M) := \max_{\lambda} \beta^{LD}(M, \lambda)$$

$$\lambda_0 = \lambda_s = 0$$

Lagrangian Relaxation provides a tighter bound than the basic lower bound, at the cost of significantly higher computational cost per node, with one exception.

# Lagrangian Relaxation for K-medoid

$$\min_{b,y} \sum_{s \in \mathcal{S}} \sum_{j \in \mathcal{S}} d_{s,j} b_{s,j}$$

$$\text{s.t.} \quad \sum_{j \in \mathcal{S}} b_{s,j} = 1$$

$$\sum_{j \in \mathcal{S}} y_j^k = 1$$

$$\sum_{k \in \mathcal{K}} y_j^k \leq 1$$

$$b_{s,j} \leq \sum_{k \in \mathcal{K}} y_j^k$$

$$b_{s,j}, y_j^k \in \{0,1\}$$

Lagrangian Relaxation

$$\beta_{LD}(M, \lambda) = \min_{b,y} \sum_{s \in \mathcal{S}} [\sum_{j \in \mathcal{S}} (d_{s,j} - \lambda_s) b_{s,j} + \lambda_s]$$

Closed-form solution (for a given  $\lambda$ ):

$$\rho_j(\lambda) := \sum_{s \in \mathcal{S}} \min(0, d_{s,j} - \lambda_s),$$

$$\beta_{LD}(M, \lambda) = \min_y \sum_{k \in \mathcal{K}} \sum_{j \in \mathcal{S}} \rho_j(\lambda) y_j^k + \sum_{s \in \mathcal{S}} \lambda_s$$

Quality guarantee of LB:

$$Gap = \frac{z - \beta_{LD}}{z_r - \beta_{LD}} < \frac{1}{e}$$

# Sample Grouping

Decompose into  
 $G = |G| \leq n$  subproblems

Non-anticipativity constraint on samples  
within groups are maintained

Proposition:  
 $\beta^B(M) \leq \beta^{SG}(M) \leq \mathcal{L}(M)$

$$\mathcal{L}(M) := \min_{m \in M} \sum_{g \in G} Q_g(m)$$

s. t.  $m_g = m_{g+1}$  ← non-anticipativity constraint



$$\beta^{SG}(M) := \min_{m_g \in M} \sum_{g \in G} Q_g(m_g)$$



$$\beta_g^{SG}(M) := \min_{m_g \in M} Q_g(m_g)$$



# Upper Bounding Methods

Basic Upper bound: Fix  $m$  at a candidate solution  $\hat{m} \in M$

Heuristic methods on root node

$$\alpha(M) := \sum_{s \in \mathcal{S}} Q_s(\hat{m})$$

- Decomposable and trivial parallelism
- Solutions from lower bound subproblems
- K-means: Lloyd Algorithm [Lloyd, 1982]
- K-center: Farthest First Traversal [Gonzalez, 1985]
- K-medoids: PAM [Kaufman and Rousseeuw, 1990]

*Theorem* [Hua et al. 2021]

With the **basic lower and upper bounds**, the BB procedure can *converge* by branching only on the centers  $\hat{m}$ .

Proof:  $Q_s(m) = \min_{k \in \mathcal{K}} \|x_s - m_k\|_2^2$  is a continuous function on  $m$ .

# Sample Determination

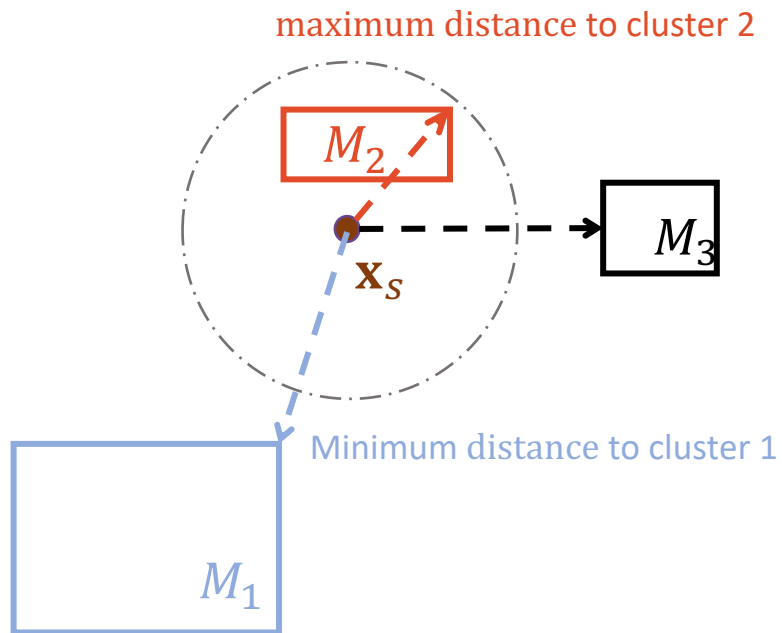
For each sample, if the following condition is satisfied, then the cluster of sample  $s$  can be determined and removed from the calculation.

$$\max_{m_{k^*} \in M_{k^*}} \|x_s - m_{k^*}\|^2 < \min_{m_k \in M_k} \|x_s - m_k\|^2, \forall k \in \mathcal{K} \setminus k^*$$

Simplify Calculation

Used for Bound Tightening

Inherited to Children Nodes



# Numerical Results for K-means [ Hua et al. 2021]

$$k\text{-means} : \min_m \sum_{S \in \mathcal{S}} \min_{k \in \mathcal{K}} \|x_S - m_k\|_2^2$$

Heuristic Benchmark: Lloyd Algorithm [Lloyd, 1982]

DATASET	SAMPLE	DIMENSION	CORES	HEURISTIC	UB	NODES	GAP(%)	TIME(S)
IRIS	150	4	20	78.85	78.85	31	$\leq 0.10$	2,160
GLASS	214	9	20	819.63	819.63	415	$\leq 0.10$	3,600
SEEDS	310	7	20	587.32	587.32	89	$\leq 0.10$	3,600
HEMI	1,955	7	20	9.75E+06	9.75E+06	112	2.23	14,400
PR2392	2,392	2	20	2.97E+10	2.97E+10	247	$\leq 0.10$	3,600
SYNTHETIC-1	42,000	2	20	5.01E+05	5.01E+05	20	3.16	14,400
SYNTHETIC-2	210,000	2	200	2.43E+06	2.43E+06	6	2.55	14,400

Practical **global optimality**  
within reasonable time

Average improvement of  
global optimal value:  
**<0.01%** (10 datasets)

Solving ML problem to global optimum is possible

Solving to global optimal solution renders a  
much better model



# Numerical Results for K-medoids [Ren et al. 2022]

$$k\text{-medoids: } \min_{m \in X} \sum_{s \in S} \min_{k \in K} \|x_s - m_k\|_2^2$$

Heuristic Benchmark: PAM [Kaufman et al., 1990]

DATASET	SAMPLE	DIMENSION	CORES	HEURISTIC	UB	NODES	GAP(%)	TIME(S)
IRIS	150	4	1	84.63	<b>83.91</b>	25	$\leq 0.10$	93
PR2392	2,392	2	1	2.13E+10	2.13E+10	37	$\leq 0.10$	123
URBANGB_10	100,000	2	40	1.26E+05	<b>1.15E+05</b>	49	$\leq 0.10$	264
RNG_AGR	199,843	7	1600	8.23E+14	8.23E+14	99	$\leq 0.10$	341
SPNET3D	434,874	3	1600	1.23E+08	2.28E+07	115	$\leq 0.10$	865
RETAIL-II	1,046,910	2	6000	2.31E+10	2.31E+10	214	$\leq 0.10$	2,515
USC1990	2,458,285	68	3000	6.91E+08	6.91E+08	25	6.33	14,400

Practical **global optimality**  
within reasonable time

Average improvement of  
global optimal value: 0.46%  
(28 datasets)

Solving ML problem to global optimum is possible

Solving to global optimal solution renders a  
much better model



# Numerical Results for K-centers [ Shi et al. 2022]

$$k\text{-centers: } \min_{m \in X} \max_{s \in S} \min_{k \in \mathcal{K}} \|x_s - m_k\|_2^2$$

Heuristic Benchmark:  
Farthest First Traversal [Gonzalez, 1985]

DATASET	SAMPLE	DIMENSION	CORES	HEURISTIC	UB	NODES	GAP(%)	TIME(S)
USC1990	2,458,285	68	1	2.04E+11	<b>1.68E+11</b>	1	$\leq 0.10$	277
HIGGS	11,000,000	29	400	368.35	<b>227.91</b>	12,576	30.2	14,400
BIGCROSS	11,620,300	56	400	1.43E+07	<b>9.38E+06</b>	10,071	$\leq 0.10$	6,444
TAXI	1,120,841,769	12	2000	3.09E+04	<b>1.62E+04</b>	1,063	$\leq 0.10$	5,705

Practical **global optimality**  
within reasonable time

Cost reduction:  
**25.8%** on average  
(38 datasets)

Solving ML problem to global optimum is possible

Solving to global optimal solution renders a  
much better model



# Reduced Space Branch and Bound

Lower Bound Strategy

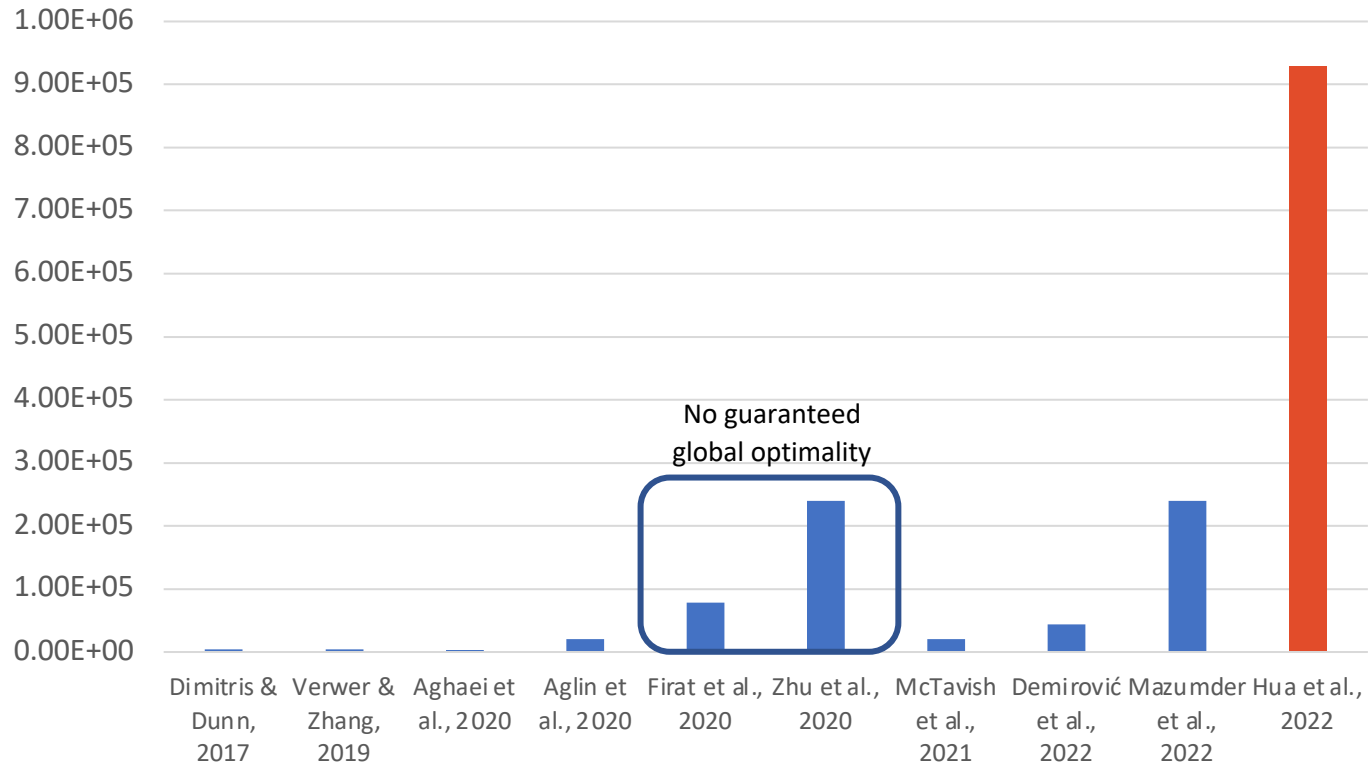
Upper Bound Strategy

Bound Tightening

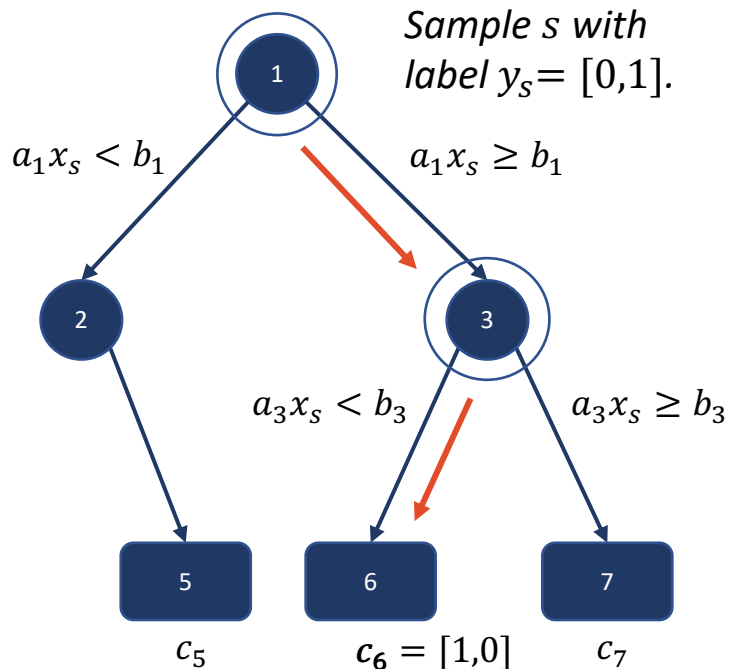
Centroid-based  
Clustering

Decision Tree

# Decision Tree Training



# Decision Tree Training



- $d_1 = d_3 = 1, d_2 = 0$ .
- $z_{s6} = 1$ .
- Since  $c_{k6} \neq y_{sk}$ , then  $L_s = 1$ .

$a, b, c, d$  describe the tree structure:

- $d_t$  where  $t \in \mathcal{T}_D$  determines whether a decision node splits or not.
- variable  $a_t = [a_{1t}, \dots, a_{pt}] \in \{0,1\}^P$  and  $b_t \in [0,1]$  are used to track the split of decision nodes.
- The prediction of each leaf node is controlled by the class indicator  $c_t \in \{0,1\}^K, \forall t \in \mathcal{T}_L$ .

$z_{st} \in \{0,1\}, \forall t \in \mathcal{T}_L, s \in \mathcal{S}$  represents whether sample  $s$  falls into leaf  $t$ .

$L_s \in [0,1]$  represents the loss of sample  $s$ .



# MILP Formulation

$$\min_{a,b,c,d,z,L} \left( \frac{1}{\hat{L}} L_s + \frac{\lambda}{n} \sum_{t \in \mathcal{T}_D} d_t \right)$$

$$\text{s.t. } \frac{1}{2} \sum_{k \in \mathcal{K}} (y_{sk} + c_{kt} - 2y_{sk}c_{kt}) - L_s \leq 1 - z_{st}, \forall t \in \mathcal{T}_L$$

$$\sum_{k \in \mathcal{K}} c_{kt} = 1, \forall t \in \mathcal{T}_L$$

$$\sum_{t \in \mathcal{T}_L} z_{st} = 1$$

$$\mathbf{a}_m^T (\mathbf{x}_s + \epsilon - \epsilon_{min}) + \epsilon_{min} \leq b_m + (1 + \epsilon_{max})(1 - z_{st}), \forall m \in A_L(t), t \in \mathcal{T}_L$$

$$\mathbf{a}_m^T \mathbf{x}_s \geq b_m - (1 - z_{st}), \forall m \in A_R(t), t \in \mathcal{T}_L$$

$$\sum_{j=1}^P a_{jt} = d_t, \forall t \in \mathcal{T}_D, j \in \{1, \dots, P\}$$

$$0 \leq b_t \leq d_t, \forall t \in \mathcal{T}_D$$

$$d_t \leq d_{p(t)}, \forall t \in \mathcal{T}_D$$

$$0 \leq L_s \leq 1$$

$$a_{jt}, d_t \in \{0, 1\}, 0 \leq b_t \leq 1 \forall t \in \mathcal{T}_D, j \in \{1, \dots, P\}$$

$$z_{st}, c_{kt} \in \{0, 1\}, \forall t \in \mathcal{T}_L$$

$$s \in \mathcal{S}$$

- Minimize the misclassification error

# MILP Formulation

$$\min_{a,b,c,d,z,L} \sum_{s \in \mathcal{S}} \left( \frac{1}{\hat{L}} L_s + \frac{\lambda}{n} \sum_{t \in \mathcal{T}_D} d_t \right)$$

$$\text{s.t. } \frac{1}{2} \sum_{k \in \mathcal{K}} (y_{sk} + c_{kt} - 2y_{sk}c_{kt}) - L_s \leq 1 - z_{st}, \forall t \in \mathcal{T}_L$$

$$\sum_{k \in \mathcal{K}} c_{kt} = 1, \forall t \in \mathcal{T}_L$$

$$\sum_{t \in \mathcal{T}_L} z_{st} = 1$$

$$\mathbf{a}_m^T (\mathbf{x}_s + \epsilon - \epsilon_{min}) + \epsilon_{min} \leq b_m + (1 + \epsilon_{max})(1 - z_{st}), \forall m \in A_L(t), t \in \mathcal{T}_L$$

$$\mathbf{a}_m^T \mathbf{x}_s \geq b_m - (1 - z_{st}), \forall m \in A_R(t), t \in \mathcal{T}_L$$

$$\sum_{j=1}^P a_{jt} = d_t, \forall t \in \mathcal{T}_D, j \in \{1, \dots, P\}$$

$$0 \leq b_t \leq d_t, \forall t \in \mathcal{T}_D$$

$$d_t \leq d_{p(t)}, \forall t \in \mathcal{T}_D$$

$$0 \leq L_s \leq 1$$

$$a_{jt}, d_t \in \{0, 1\}, 0 \leq b_t \leq 1 \quad \forall t \in \mathcal{T}_D, j \in \{1, \dots, P\}$$

$$z_{st}, c_{kt} \in \{0, 1\}, \forall t \in \mathcal{T}_L$$

$$s \in \mathcal{S}$$

First Stage Variable

$m : a, b, c, d$

(Tree structure Variables, mixed integer)

Restrict each sample to be predicted  
follows a decision tree structure

Second Stage Variable

$z_{s,t}, L_s$

(Leaf assignment and loss, mixed integer)

# Basic Lower Bound

$$\beta_S(M) := \min_{m_s \in M} Q_S(m_s)$$
$$m := (a, b, c, d)$$

Decompose into  $|\mathcal{S}| =$   
 $n$  subproblems



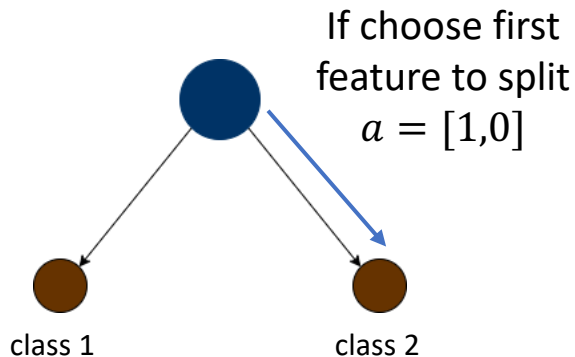
Searching for  
reachable leaves

# Basic Lower Bound

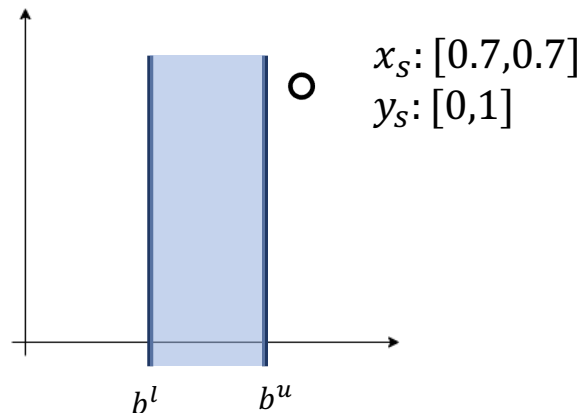
Given a branching node (in the branch and bound process) with bound  $M$ :

$$a^l = [0,0], \quad b^l = [0.3], \quad c^l = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad d^l = [1]$$

$$a^u = [1,1], \quad b^u = [0.6], \quad c^u = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad d^u = [1]$$



$$\mathcal{J}_{z_s} = \{3\}, \beta_s = 0$$

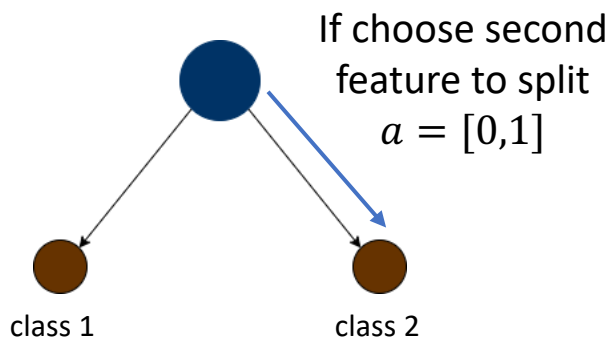


# Basic Lower Bound

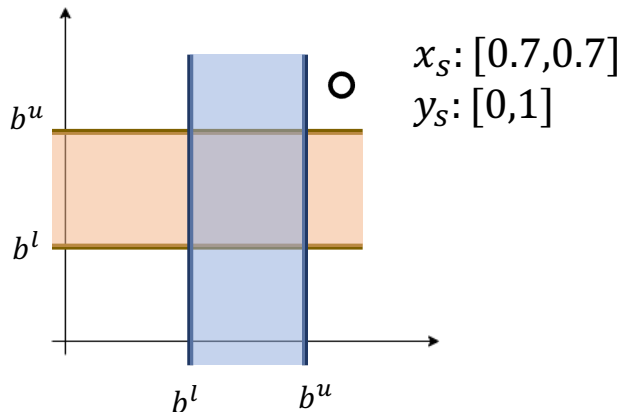
Given a branching node (in the branch and bound process) with bound  $M$ :

$$a^l = [0,0], \quad b^l = [0.3], \quad c^l = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad d^l = [1]$$

$$a^u = [1,1], \quad b^u = [0.6], \quad c^u = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad d^u = [1]$$



$$\mathcal{J}_{z_s} = \{3\}, \beta_s = 0$$

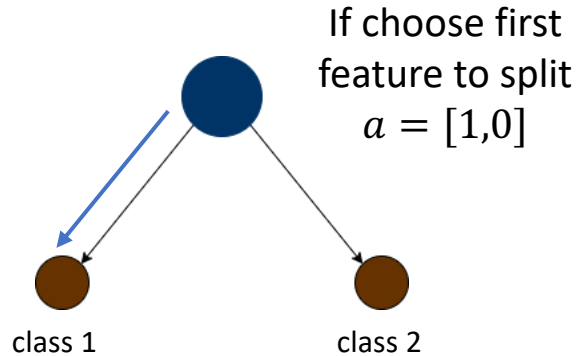


# Basic Lower Bound

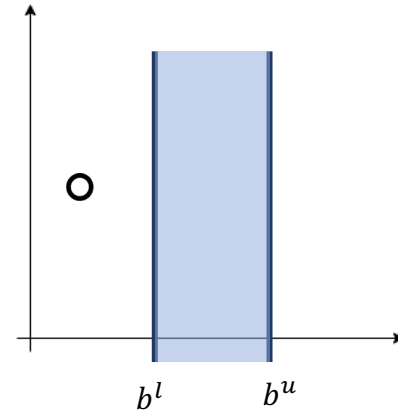
Given a branching node (in the branch and bound process) with bound  $M$ :

$$a^l = [0,0], \quad b^l = [0.3], \quad c^l = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad d^l = [1]$$

$$a^u = [1,1], \quad b^u = [0.6], \quad c^u = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad d^u = [1]$$



$$x_s: [0.2, 0.5]$$
$$y_s: [0, 1]$$

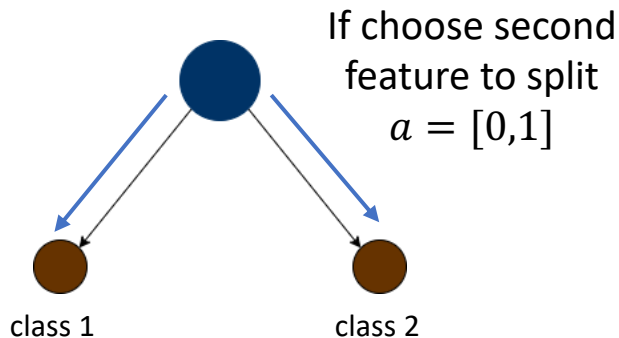


# Basic Lower Bound

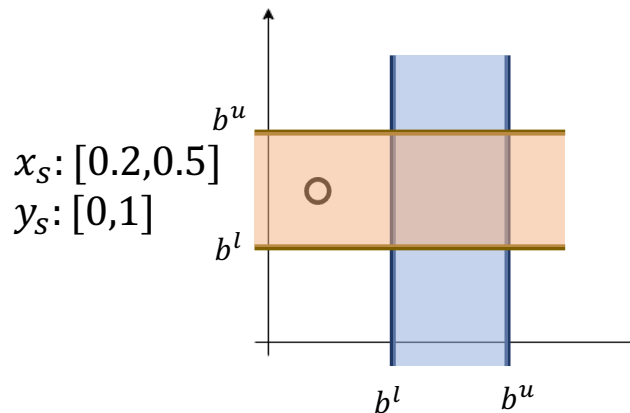
Given a branching node (in the branch and bound process) with bound  $M$ :

$$a^l = [0,0], \quad b^l = [0.3], \quad c^l = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad d^l = [1]$$

$$a^u = [1,1], \quad b^u = [0.6], \quad c^u = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad d^u = [1]$$



$$\mathcal{J}_{z_s} = \{2,3\}, \beta_s = 0$$



# Sample Determination

- Comparing the label  $y_{sk}$  with the range of  $c$  of each leaf node in  $t \in \mathcal{T}_{z_s}$ , the loss of some samples can be determined.
- Suppose  $k$  is the true label of sample  $s$  (i.e.,  $y_{sk} = 1$ ), then we can check the loss of sample  $s$  with the following equation:

$$L_s = \begin{cases} 1 & \text{if } \bigwedge_{t \in \mathcal{T}_{z_s}} c_{kt}^l = c_{kt}^u = 0 \\ 0 & \text{if } \bigwedge_{t \in \mathcal{T}_{z_s}} c_{kt}^l = c_{kt}^u = 1 \\ \text{undtm.} & \text{otherwise} \end{cases}$$



# Numerical Results (D=2, nc=1000) [Hua et al. 2022]

Benchmark:

CART [Breiman et al, 1984],  
OCT[Bertimas and Dunn, 2019]

Practical **global optimality**  
within reasonable time

DATA-SET	S	P	K	METHOD	UB	TRAINING ACCURACY(%)	TESTING ACCURACY(%)	GAP (%)	TIME (s)
PENDIGITS	7,494	16	10	CART	32993.0	38.9	39.1	-	-
				OCT	32993.0	38.9	39.1	100	14400
				RS-OCT	<b>32426.1</b>	<b>40.0</b>	<b>39.9</b>	<b>&lt;1%</b>	<b>2093.5</b>
AVILA	10,430	10	12	CART	9454.3	50.4	51.0	-	-
				OCT	9196.4	51.7	52.5	100	14400
				RS-OCT	<b>8787.5</b>	<b>53.8</b>	<b>53.8</b>	<b>19.6</b>	14400
EEG	14,980	14	2	CART	7815.8	61.7	61.2	-	-
				OCT	7748.7	61.7	61.0	100	14400
				RS-OCT	<b>6752.7</b>	<b>66.9</b>	<b>65.5</b>	<b>50.8</b>	14400
HTRU	17,898	8	2	CART	327.0	97.8	97.7	-	-
				OCT	320.5	97.8	97.7	100	14400
				RS-OCT	<b>300.7</b>	<b>98.0</b>	<b>97.8</b>	<b>42.8</b>	14400
SHUTTLE	43,500	9	7	CART	2567.5	93.8	93.8	-	-
				OCT	2567.5	93.8	93.8	100	14400
				RS-OCT	<b>1908.1</b>	<b>95.4</b>	<b>95.5</b>	<b>&lt;1%</b>	<b>586.1</b>
SKIN-SEGMENTATION	245,057	3	2	CART	23409.5	89.9	89.8	-	-
				OCT	23409.5	89.9	89.8	100	14400
				RS-OCT	<b>16953.7</b>	<b>92.7</b>	<b>92.7</b>	<b>&lt;1%</b>	<b>6698.9</b>
HT-SENSOR	928,991	11	3	CART	782046.1	58.1	58.1	-	-
				OCT	782046.1	58.1	58.1	100	14400
				RS-OCT	<b>753075.8</b>	<b>59.7</b>	<b>59.7</b>	<b>56.2</b>	14400

Solving ML problem to global optimum is possible



# Numerical Results (D=2, nc=1000) [Hua et al. 2022]

## Benchmark:

CART [Breiman et al, 1984],  
OCT[Bertimas and Dunn, 2019]

Practical **global optimality**  
within reasonable time

Cost reduction:  
**12.5%** and **11.7%**  
on average

Solving ML problem to global  
optimum is possible



Solving to global optimal solution renders a  
much better model



DATA-SET	S	P	K	METHOD	UB	TRAINING ACCURACY(%)	TESTING ACCURACY(%)	GAP (%)	TIME (s)
PENDIGITS	7,494	16	10	CART	32993.0	38.9	39.1	-	-
				OCT	32993.0	38.9	39.1	100	14400
				RS-OCT	<b>32426.1</b>	<b>40.0</b>	<b>39.9</b>	<b>&lt;1%</b>	<b>2093.5</b>
AVILA	10,430	10	12	CART	9454.3	50.4	51.0	-	-
				OCT	9196.4	51.7	52.5	100	14400
				RS-OCT	<b>8787.5</b>	<b>53.8</b>	<b>53.8</b>	<b>19.6</b>	14400
EEG	14,980	14	2	CART	7815.8	61.7	61.2	-	-
				OCT	7748.7	61.7	61.0	100	14400
				RS-OCT	<b>6752.7</b>	<b>66.9</b>	<b>65.5</b>	<b>50.8</b>	14400
HTRU	17,898	8	2	CART	327.0	97.8	97.7	-	-
				OCT	320.5	97.8	97.7	100	14400
				RS-OCT	<b>300.7</b>	<b>98.0</b>	<b>97.8</b>	<b>42.8</b>	14400
SHUTTLE	43,500	9	7	CART	2567.5	93.8	93.8	-	-
				OCT	2567.5	93.8	93.8	100	14400
				RS-OCT	<b>1908.1</b>	<b>95.4</b>	<b>95.5</b>	<b>&lt;1%</b>	<b>586.1</b>
SKIN-SEGMENTATION	245,057	3	2	CART	23409.5	89.9	89.8	-	-
				OCT	23409.5	89.9	89.8	100	14400
				RS-OCT	<b>16953.7</b>	<b>92.7</b>	<b>92.7</b>	<b>&lt;1%</b>	<b>6698.9</b>
HT-SENSOR	928,991	11	3	CART	782046.1	58.1	58.1	-	-
				OCT	782046.1	58.1	58.1	100	14400
				RS-OCT	<b>753075.8</b>	<b>59.7</b>	<b>59.7</b>	<b>56.2</b>	14400

# Numerical Results (D=2, nc=1000) [Hua et al. 2022]

Benchmark:

CART [Breiman et al, 1984],  
OCT[Bertimas and Dunn, 2019]

Training accuracy:

**3.7%** and **3.6%** on average

Testing accuracy:

**3.3%** and **2.8%** on average

DATA-SET	S	P	K	METHOD	UB	TRAINING ACCURACY(%)	TESTING ACCURACY(%)	GAP (%)	TIME (s)
PENDIGITS	7,494	16	10	CART	32993.0	38.9	39.1	-	-
				OCT	32993.0	38.9	39.1	100	14400
				RS-OCT	<b>32426.1</b>	<b>40.0</b>	<b>39.9</b>	<b>&lt;1%</b>	<b>2093.5</b>
AVILA	10,430	10	12	CART	9454.3	50.4	51.0	-	-
				OCT	9196.4	51.7	52.5	100	14400
				RS-OCT	<b>8787.5</b>	<b>53.8</b>	<b>53.8</b>	<b>19.6</b>	14400
EEG	14,980	14	2	CART	7815.8	61.7	61.2	-	-
				OCT	7748.7	61.7	61.0	100	14400
				RS-OCT	<b>6752.7</b>	<b>66.9</b>	<b>65.5</b>	<b>50.8</b>	14400
HTRU	17,898	8	2	CART	327.0	97.8	97.7	-	-
				OCT	320.5	97.8	97.7	100	14400
				RS-OCT	<b>300.7</b>	<b>98.0</b>	<b>97.8</b>	<b>42.8</b>	14400
SHUTTLE	43,500	9	7	CART	2567.5	93.8	93.8	-	-
				OCT	2567.5	93.8	93.8	100	14400
				RS-OCT	<b>1908.1</b>	<b>95.4</b>	<b>95.5</b>	<b>&lt;1%</b>	<b>586.1</b>
SKIN-SEGMENTATION	245,057	3	2	CART	23409.5	89.9	89.8	-	-
				OCT	23409.5	89.9	89.8	100	14400
				RS-OCT	<b>16953.7</b>	<b>92.7</b>	<b>92.7</b>	<b>&lt;1%</b>	<b>6698.9</b>
HT-SENSOR	928,991	11	3	CART	782046.1	58.1	58.1	-	-
				OCT	782046.1	58.1	58.1	100	14400
				RS-OCT	<b>753075.8</b>	<b>59.7</b>	<b>59.7</b>	<b>56.2</b>	14400

Solving ML problem to global optimum is possible



Solving to global optimal solution renders a much better model



# Numerical Results of DE-MH (D=2)

Benchmark:  
RS-OCT [Hua et al. 2022]

Training accuracy:  
DE-MH (GPU) is **0.10% smaller**  
than RS-OCT on average

Testing accuracy:  
DE-MH (GPU) is **0.04% larger**  
than RS-OCT on average

For datasets with  $\leq 1\%$  gap,  
**0.25% and 0.39% smaller** on  
training and testing accuracy

DATASET	SAMPLE	DIMEN- SION	CLASS	METHOD	TRAINING ACCURACY(%)	TESTING ACCURACY(%)	GAP (%)	TIME (S)
PEN- DIGITS	7,494	16	10	RS-OCT	<b>40.0</b>	<b>39.9</b>	$\leq 1\%$	2,094
				DE-MH (CPU)	<b>40.0</b>	39.8	-	72
				DE-MH (GPU)	<b>40.0</b>	39.8	-	13
AVILA	10,430	10	12	RS-OCT	<b>53.8</b>	53.8	19.6	14,400
				DE-MH (CPU)	53.7	<b>53.9</b>	-	79
				DE-MH (GPU)	53.7	53.7	-	10
EEG	14,980	14	2	RS-OCT	66.9	65.5	50.8	14,400
				DE-MH (CPU)	66.5	65.8	-	101
				DE-MH (GPU)	<b>67.0</b>	<b>66.6</b>	-	10
HTRU	17,898	8	2	RS-OCT	<b>98.0</b>	<b>97.8</b>	42.8	14,400
				DE-MH (CPU)	97.9	97.7	-	94
				DE-MH (GPU)	97.9	97.6	-	10
SHUTTLE	43,500	9	7	RS-OCT	<b>95.4</b>	<b>95.5</b>	$\leq 1\%$	586
				DE-MH (CPU)	94.6	94.4	-	276
				DE-MH (GPU)	94.7	94.5	-	11
SKIN- SEGMENT- TATION	245,057	3	2	RS-OCT	<b>92.7</b>	<b>92.7</b>	$\leq 1\%$	6,699
				DE-MH (CPU)	<b>92.7</b>	<b>92.7</b>	-	973
				DE-MH (GPU)	<b>92.7</b>	<b>92.7</b>	-	12
HT- SENSOR	928,991	11	3	RS-OCT	59.7	59.7	56.2	14,400
				DE-MH (CPU)	<b>59.8</b>	<b>59.8</b>	-	6,074
				DE-MH (GPU)	<b>59.8</b>	<b>59.8</b>	-	32

DE-MH (GPU) can maintain  
**comparable levels of accuracy** to the  
global method RS-OCT



# Numerical Results of DE-MH (D=2)

Benchmark:  
RS-OCT [Hua et al. 2022]

Running Time:  
DE-MH (GPU) offers a **675X**  
improvement in speed over  
RS-OCT with 1000 cores

DATASET	SAMPLE	DIMEN- SION	CLASS	METHOD	TRAINING ACCURACY(%)	TESTING ACCURACY(%)	GAP (%)	TIME (S)
PEN- DIGITS	7,494	16	10	RS-OCT	<b>40.0</b>	<b>39.9</b>	≤1%	2,094
				DE-MH (CPU)	<b>40.0</b>	39.8	-	72
				DE-MH (GPU)	<b>40.0</b>	39.8	-	13
AVILA	10,430	10	12	RS-OCT	<b>53.8</b>	53.8	19.6	14,400
				DE-MH (CPU)	53.7	<b>53.9</b>	-	79
				DE-MH (GPU)	53.7	53.7	-	10
EEG	14,980	14	2	RS-OCT	66.9	65.5	50.8	14,400
				DE-MH (CPU)	66.5	65.8	-	101
				DE-MH (GPU)	<b>67.0</b>	<b>66.6</b>	-	10
HTRU	17,898	8	2	RS-OCT	<b>98.0</b>	<b>97.8</b>	42.8	14,400
				DE-MH (CPU)	97.9	97.7	-	94
				DE-MH (GPU)	97.9	97.6	-	10
SHUTTLE	43,500	9	7	RS-OCT	<b>95.4</b>	<b>95.5</b>	≤1%	586
				DE-MH (CPU)	94.6	94.4	-	276
				DE-MH (GPU)	94.7	94.5	-	11
SKIN- SEGMENTATION	245,057	3	2	RS-OCT	<b>92.7</b>	<b>92.7</b>	≤1%	6,699
				DE-MH (CPU)	<b>92.7</b>	<b>92.7</b>	-	973
				DE-MH (GPU)	<b>92.7</b>	<b>92.7</b>	-	12
HT- SENSOR	928,991	11	3	RS-OCT	59.7	59.7	56.2	14,400
				DE-MH (CPU)	<b>59.8</b>	<b>59.8</b>	-	6,074
				DE-MH (GPU)	<b>59.8</b>	<b>59.8</b>	-	32

The metaheuristic method obtains **comparable levels of accuracy** to the global method with **significantly reduced running time.**



## Numerical Results of DE-MH (D=3)

Benchmark:

CART [Breiman et al, 1984],  
RS-OCT [Hua et al. 2022]

DE-MH v.s. CART

3.36% and 3.26% larger on  
training and testing accuracy

DE-MH v.s. RS-OCT

2.52% and 2.75% larger on  
training and testing accuracy  
(RS-OCT fails to achieve reasonable  
gaps)

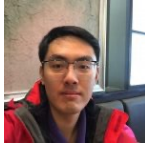
DATASET	SAMPLE	DIMEN- SION	CLASS	METHOD	TRAINING ACCURACY(%)	TESTING ACCURACY(%)	GAP (%)	TIME (S)
PEN- DIGITS	7,494	16	10	CART	62.5	63.0	-	0.01
				RS-OCT	63.9	63.7	46.7	14,400
				DE-MH (GPU)	<b>66.9</b>	<b>66.8</b>	-	32
AVILA	10,430	10	12	CART	54.4	54.6	-	0.02
				RS-OCT	55.0	55.3	91.1	14,400
				DE-MH (GPU)	<b>57.2</b>	<b>57.4</b>	-	28
EEG	14,980	14	2	CART	66.6	65.2	-	0.02
				RS-OCT	67.9	65.7	100	14,400
				DE-MH (GPU)	<b>70.4</b>	<b>69.2</b>	-	27
HTRU	17,898	8	2	CART	97.9	97.7	-	0.04
				RS-OCT	97.9	97.7	99.9	14,400
				DE-MH (GPU)	<b>98.0</b>	<b>97.8</b>	-	21
SHUTTLE	43,500	9	7	CART	99.7	99.6	-	0.03
				RS-OCT	99.7	<b>99.8</b>	84.9	14,400
				DE-MH (GPU)	<b>99.8</b>	<b>99.8</b>	-	30
SKIN- SEGMENT- TATION	245,057	3	2	CART	96.5	96.5	-	0.07
				RS-OCT	96.6	96.5	80.3	14,400
				DE-MH (GPU)	<b>96.7</b>	<b>96.8</b>	-	29
HT- SENSOR	928,991	11	3	CART	64.5	64.5	-	2.71
				RS-OCT	64.6	64.5	100	14,400
				DE-MH (GPU)	<b>67.7</b>	<b>67.7</b>	-	65

# Conclusion and Highlights

- A reduced-space branch and bound training algorithm.
  - Guaranteed convergence by **branching only on model parameter variables**.
- Decomposable lower and upper bound strategies.
  - Trivial Parallelism with closed-form subproblem solutions.
- Achieve global optimum solutions for large-scale ML problems.
  - **One billion samples** within 2 hours (0.1%, k-center problem, 1000 cores).
- Global Optimization vs Heuristics

# Acknowledge

## Postdoc



Kaixun Hua



Kevin Kung



Yan Yu



Vijay Pediredla

## Phd



Yixiu Wang



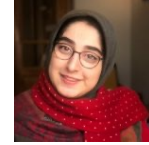
Jiayang Ren



Chaojie Ji



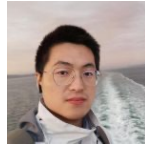
Qiangqiang Mao



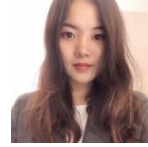
Ghazaleh Mozafari



Mohammad Aliahmadi



Liang Cao



Jingyi Wang

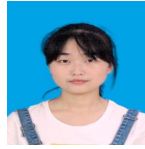


Ahmed Al-Shafei

## Master



Morimasa Okamoto



Mingfei Shi



Jeff MacDonald





**Thank You!**

