

Supervised Learning:

# **DECISION TREES (RECURSIVE PARTITIONING)**

# Decision Trees

- supervised learning
- classification or regression
- traditionally created manually
- tree-growing algorithms to optimize decision

Example: Processing mortgage applications

YRSJOB: Years at current job

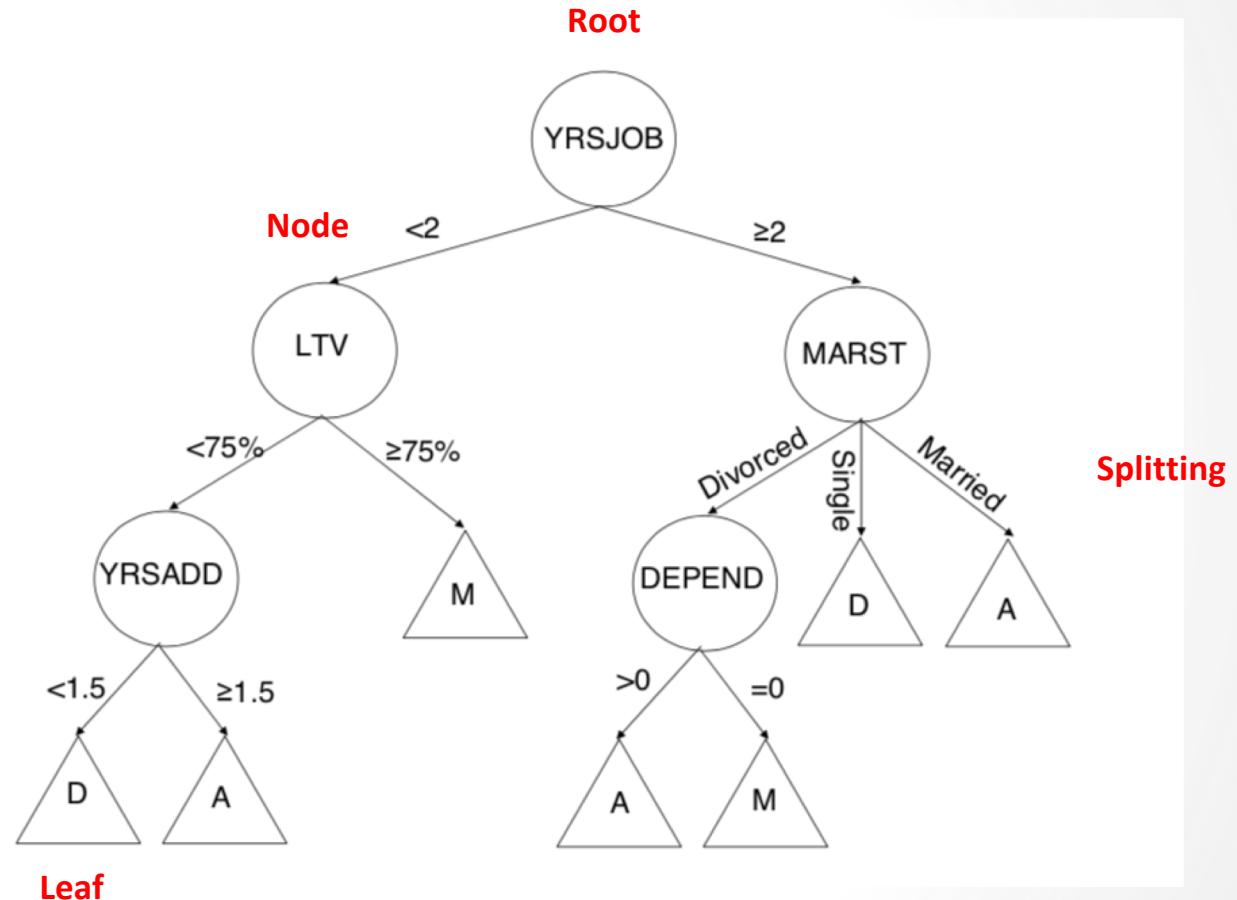
YRSADD: Years at current address

DEPEND: # of dependents

MARST: marital status

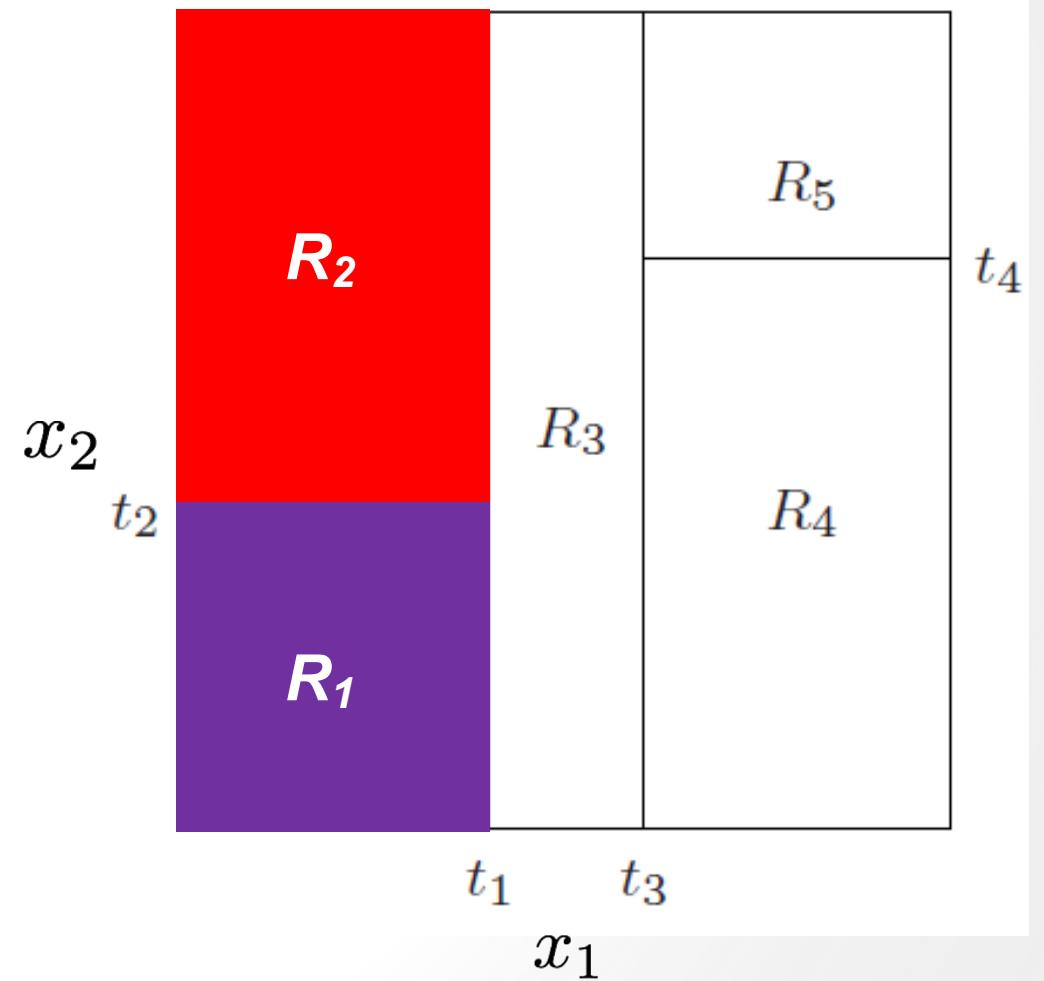
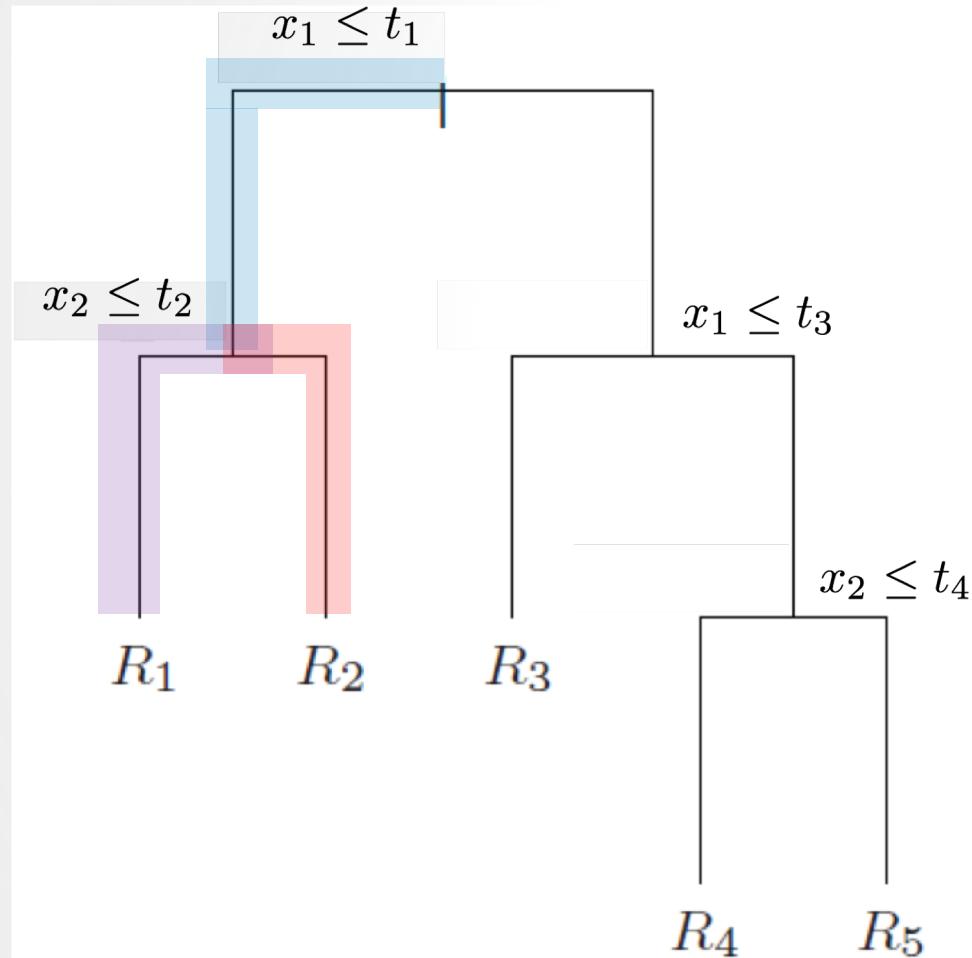
LTВ: loan-to-value ratio

$$T : \mathbf{x} \rightarrow y$$



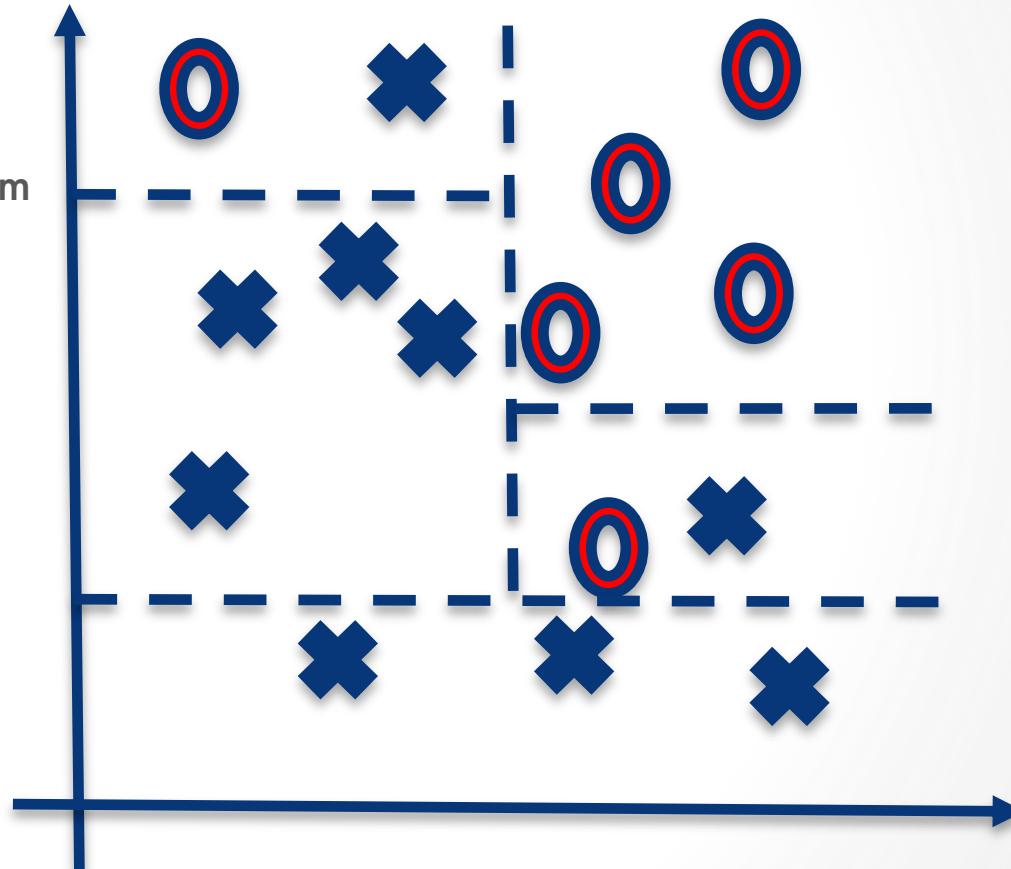
# Decision Trees

- Decision tree is a recursive partition of the input space
- Most are *binary*



# Growing a Tree: Stopping Criteria

- 1) All instances in leaf have the same y value
- 2) A maximum tree depth has been reached
- 3) Number of instances in leaf is below some minimum
- 4) Splitting criteria is below some threshold



# Missing and Unseen Values

- Missing Data in Training Predictors: Treated as a “level”
- Missing Value in Training Target: Ignored
- Unseen Categorical Levels During Scoring: Treated as an “NA” value and classified with outliers

# Pros/Cons of Decision Trees

simple to understand/interpret

little data prep

- natural handling of "mixed" data types
- handling of missing values
- handing of multi-class outputs

robustness to outliers in input space

insensitive to monotonic transformations of inputs

computational scalability (large N)  $\mathcal{O}(n \log n)$

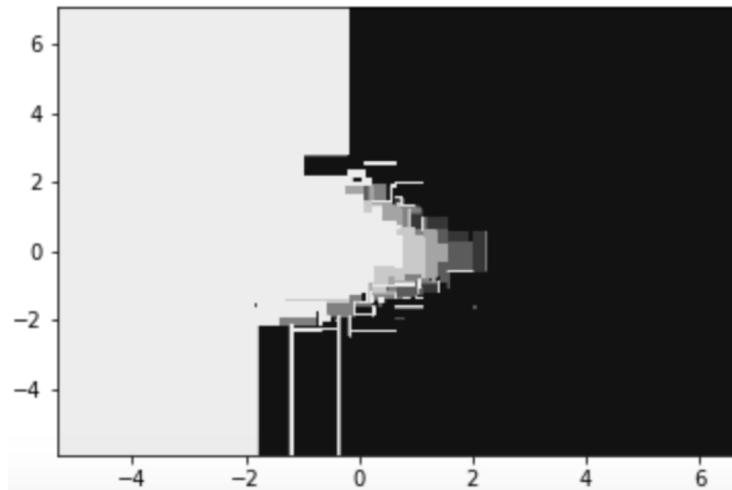
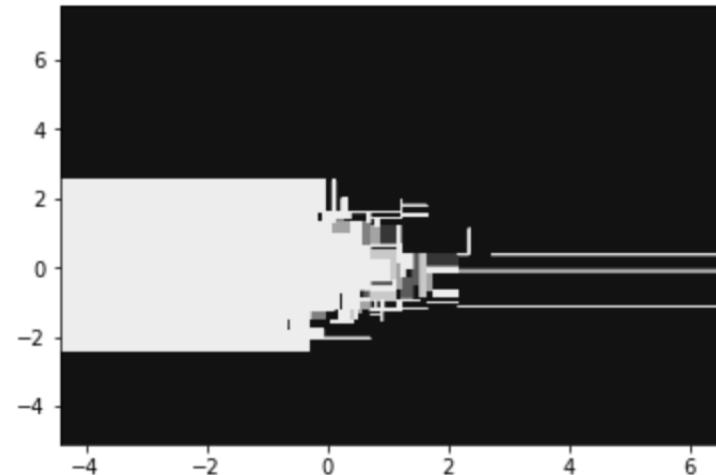
automatically ignores irrelevant inputs

**X** weak predictors

**X** can be unstable to small variations in the data

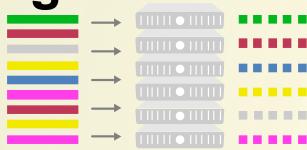
**X** poor ability to extract linear relationships

**X** can create biased trees if classes unbalanced



# Scalable Implementation in H2O

## 1 Parallel Data Ingest



Data is stored in-memory on all cluster compute nodes

- Rows are evenly distributed across the cluster
- Columns are stored separately and compressed

Basis for fine-grain Map/Reduce for histogram calculation

## 2 Distributed Tree Building

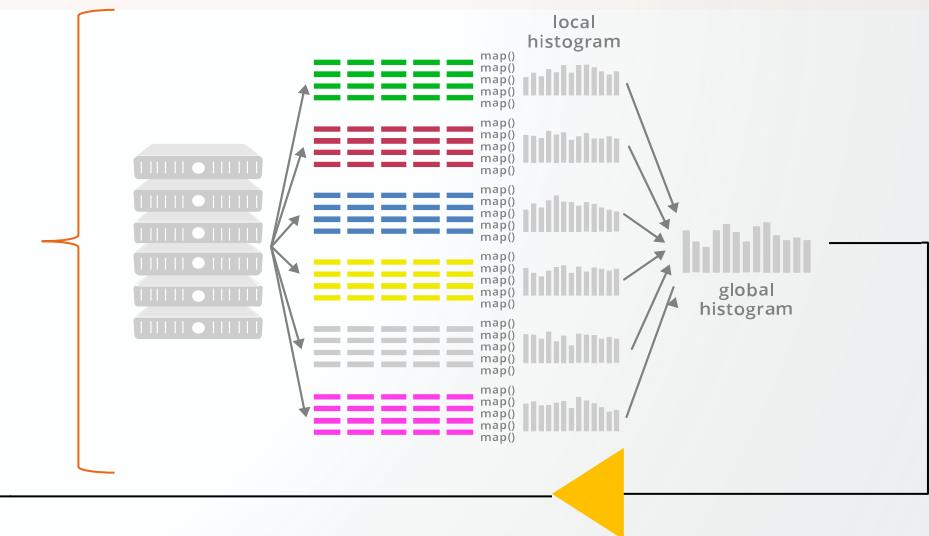
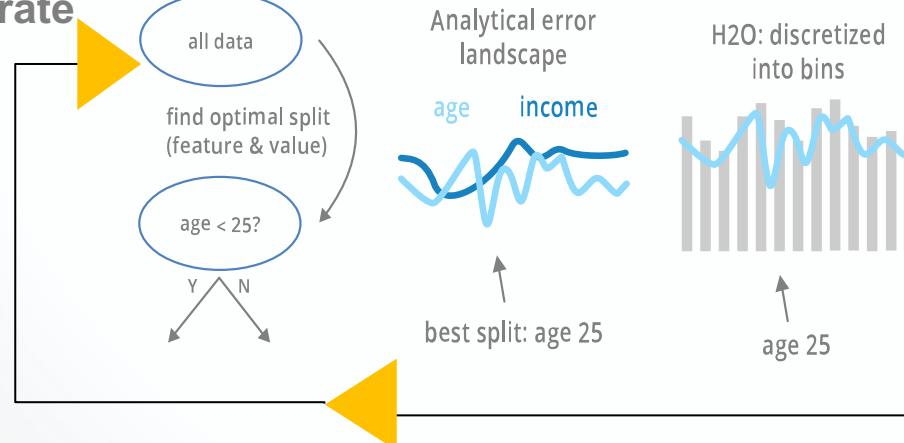
via Fine-Grain Map/Reduce to find optimal split points of data layer by layer

Start with **root node** and build layers of tree nodes [ILLUSTRATION BELOW]

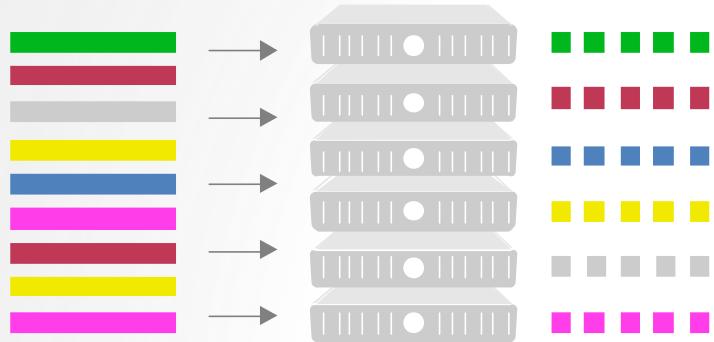
For each layer, **repeat** the following:

- For a set of features, split the data at every possible split point
- Find the split that leads to best model improvement
- Use discretization to limit the number of potential splits
  - To find the split, local histograms are calculated on each node and then aggregated into a global histogram
  - From the global histogram, the best split column is chosen

For each layer,  
iterate



# Scalable Implementation in H2O

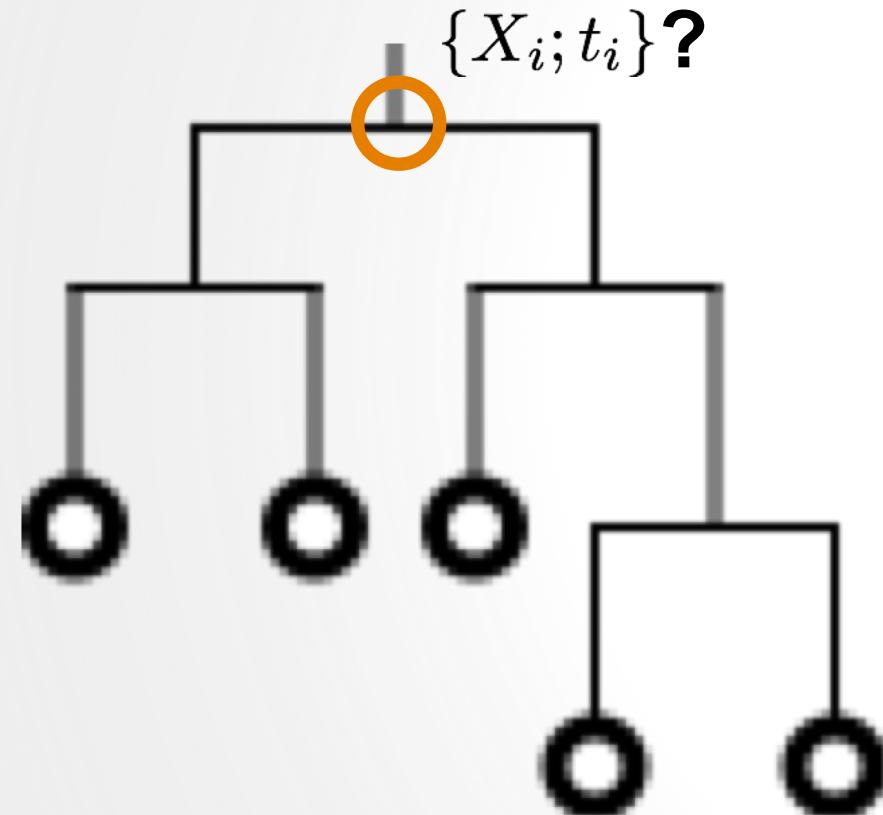


Parallel Parse into **Distributed Rows**

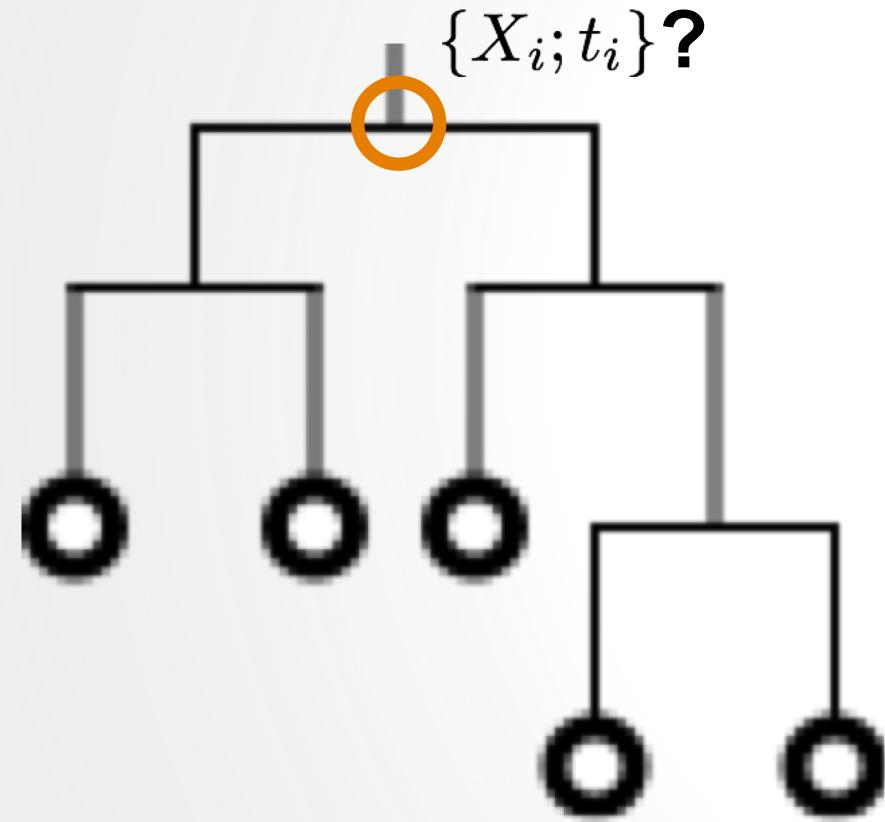


**Fine Grain Map Reduce Illustration:** Scalable Distributed Histogram Calculation for GBM

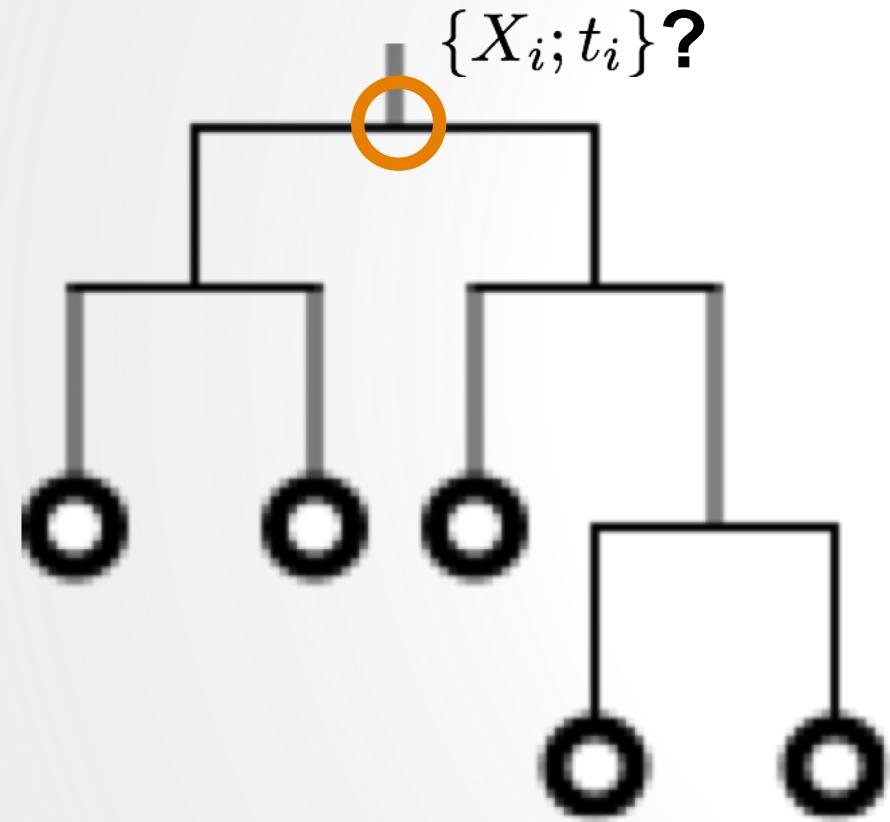
# Tree Growth with Data Parallelism



# Tree Growth with Data Parallelism

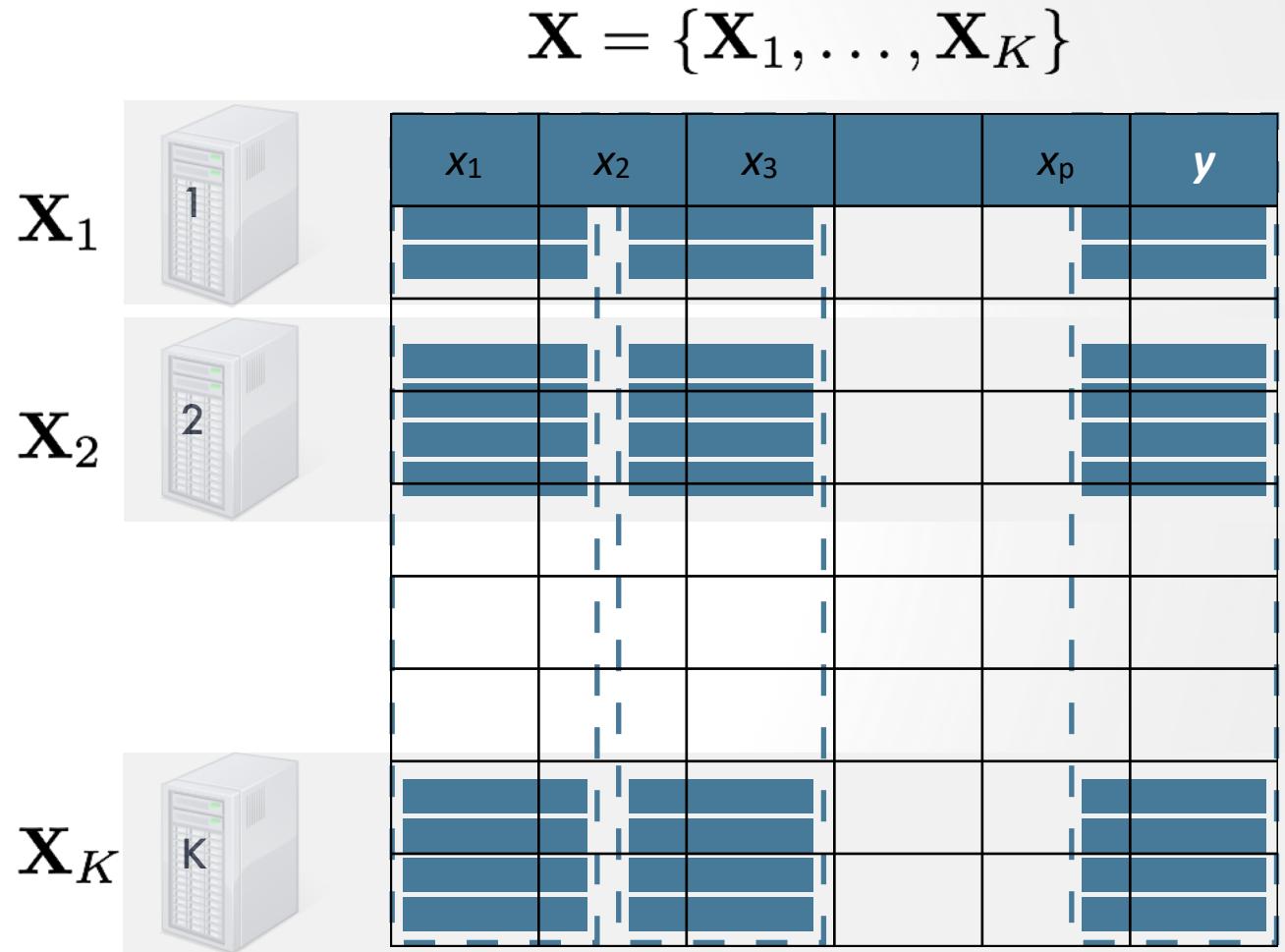
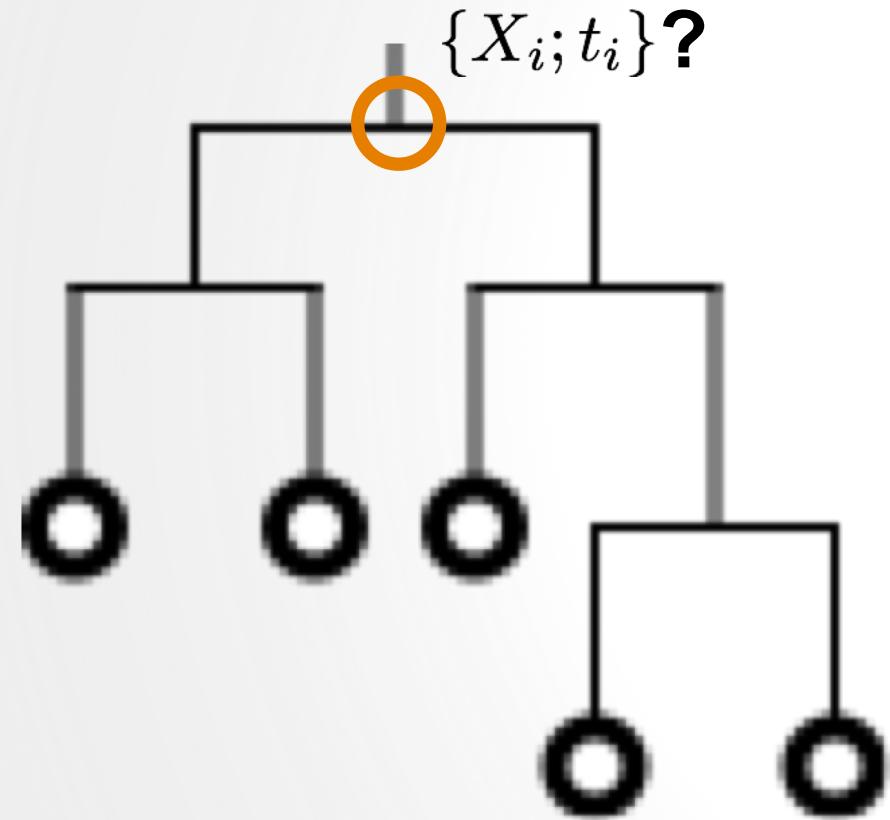


# Tree Growth with Data Parallelism

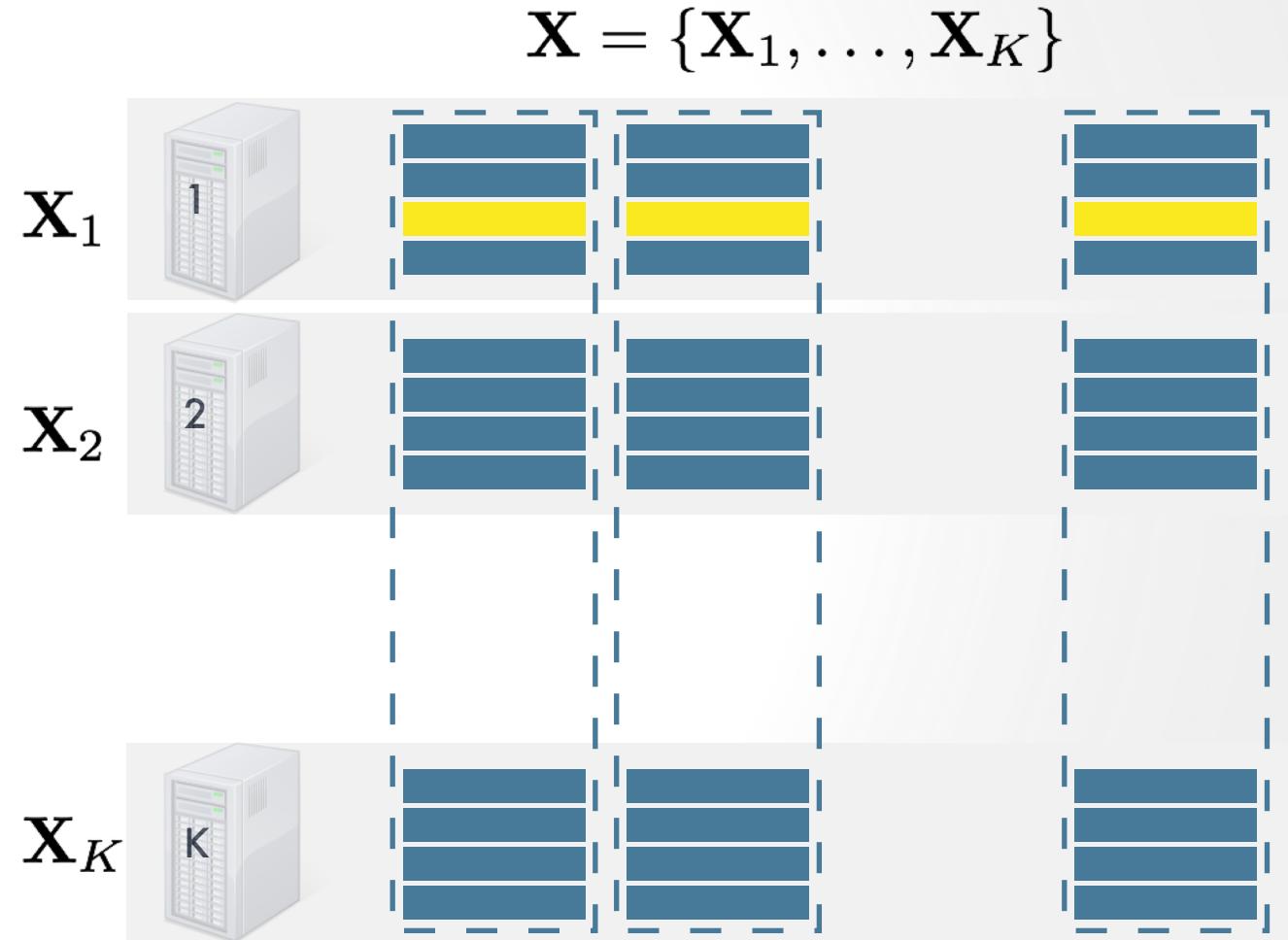
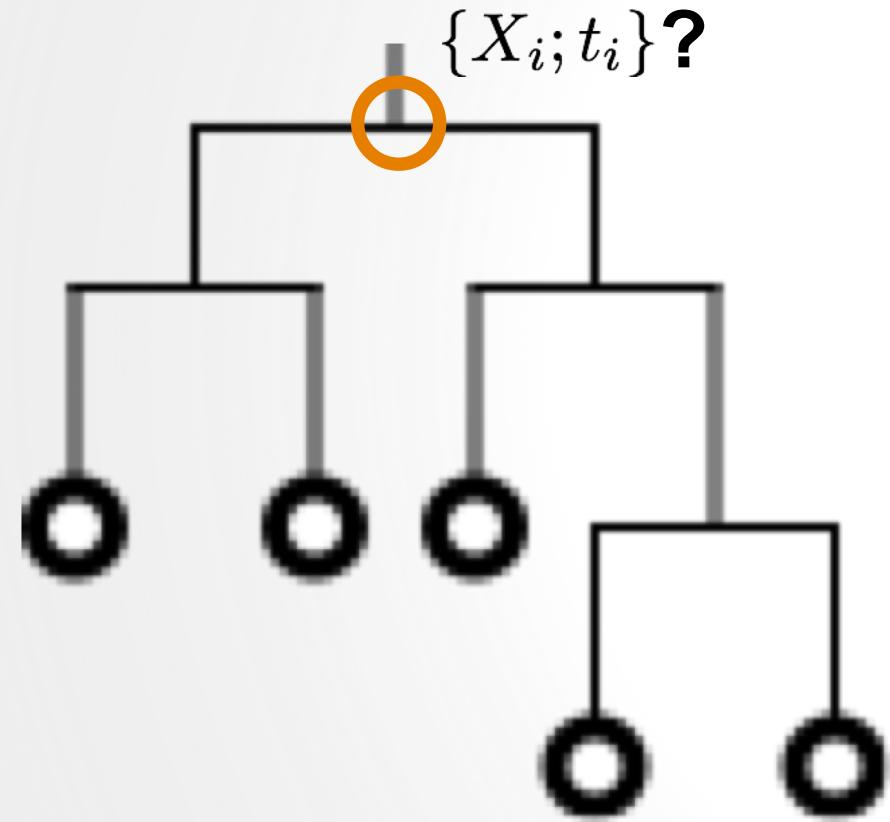


$x_1$	$x_2$	$x_3$		$x_p$	$y$
[Data Row 1]	[Data Row 1]	[Data Row 1]		[Data Row 1]	[Data Row 1]
[Data Row 2]	[Data Row 2]	[Data Row 2]		[Data Row 2]	[Data Row 2]
[Data Row 3]	[Data Row 3]	[Data Row 3]		[Data Row 3]	[Data Row 3]
[Data Row 4]	[Data Row 4]	[Data Row 4]		[Data Row 4]	[Data Row 4]
[Data Row 5]	[Data Row 5]	[Data Row 5]		[Data Row 5]	[Data Row 5]
[Data Row 6]	[Data Row 6]	[Data Row 6]		[Data Row 6]	[Data Row 6]
[Data Row 7]	[Data Row 7]	[Data Row 7]		[Data Row 7]	[Data Row 7]
[Data Row 8]	[Data Row 8]	[Data Row 8]		[Data Row 8]	[Data Row 8]
[Data Row 9]	[Data Row 9]	[Data Row 9]		[Data Row 9]	[Data Row 9]
[Data Row 10]	[Data Row 10]	[Data Row 10]		[Data Row 10]	[Data Row 10]
[Data Row 11]	[Data Row 11]	[Data Row 11]		[Data Row 11]	[Data Row 11]
[Data Row 12]	[Data Row 12]	[Data Row 12]		[Data Row 12]	[Data Row 12]
[Data Row 13]	[Data Row 13]	[Data Row 13]		[Data Row 13]	[Data Row 13]
[Data Row 14]	[Data Row 14]	[Data Row 14]		[Data Row 14]	[Data Row 14]
[Data Row 15]	[Data Row 15]	[Data Row 15]		[Data Row 15]	[Data Row 15]
[Data Row 16]	[Data Row 16]	[Data Row 16]		[Data Row 16]	[Data Row 16]
[Data Row 17]	[Data Row 17]	[Data Row 17]		[Data Row 17]	[Data Row 17]
[Data Row 18]	[Data Row 18]	[Data Row 18]		[Data Row 18]	[Data Row 18]
[Data Row 19]	[Data Row 19]	[Data Row 19]		[Data Row 19]	[Data Row 19]
[Data Row 20]	[Data Row 20]	[Data Row 20]		[Data Row 20]	[Data Row 20]

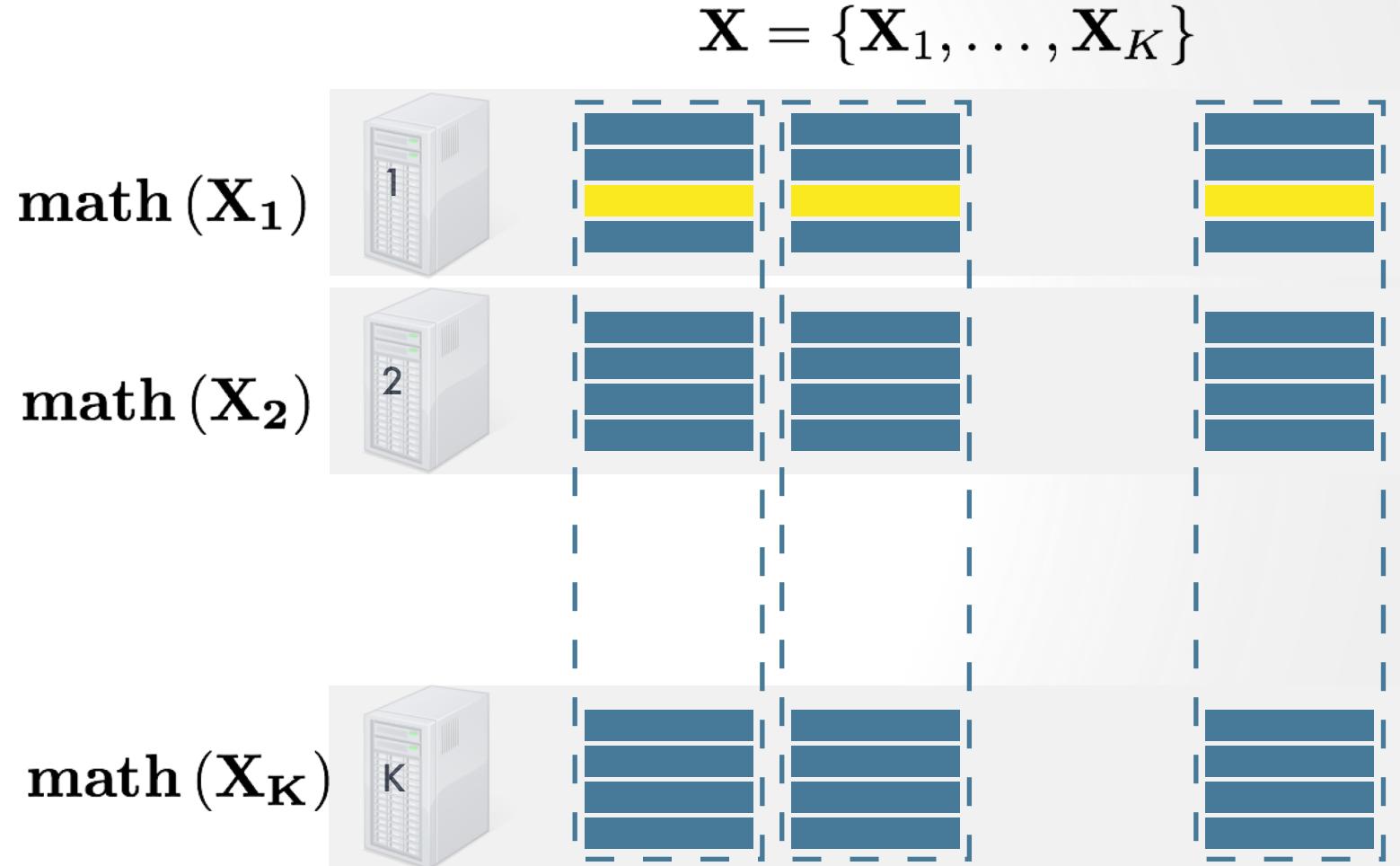
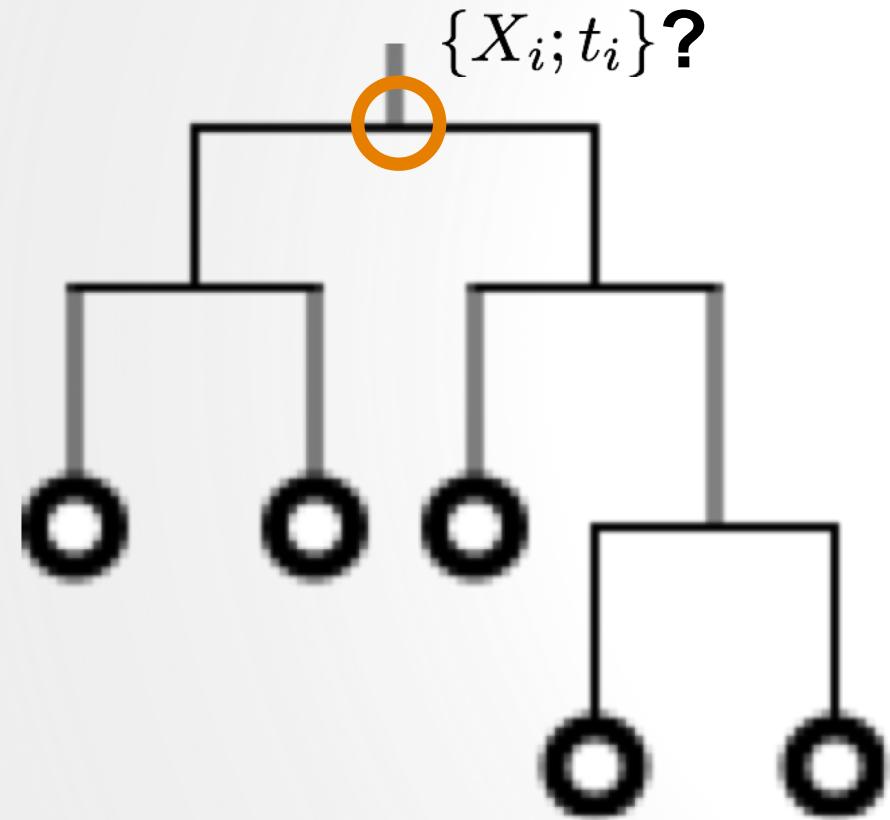
# Tree Growth with Data Parallelism



# Tree Growth with Data Parallelism

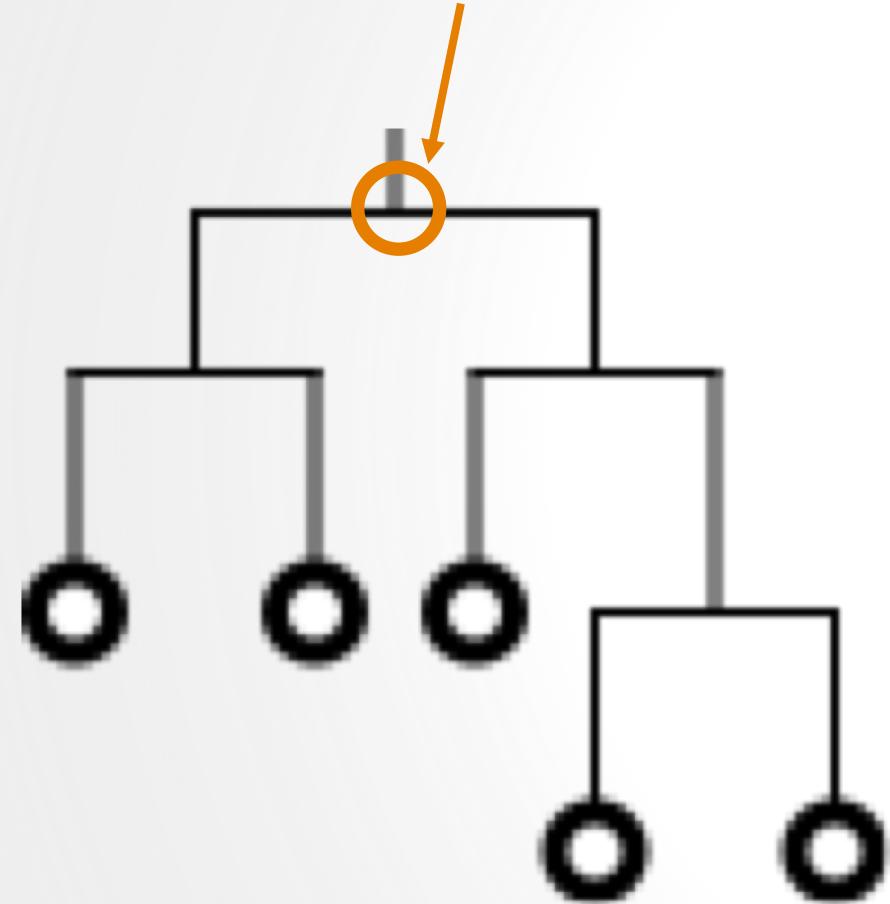


# Tree Growth with Data Parallelism



# Tree Growth with Data Parallelism

$$\{X_i; t_i\} = f(\mathbf{math}(X_1), \dots, \mathbf{math}(X_K))$$

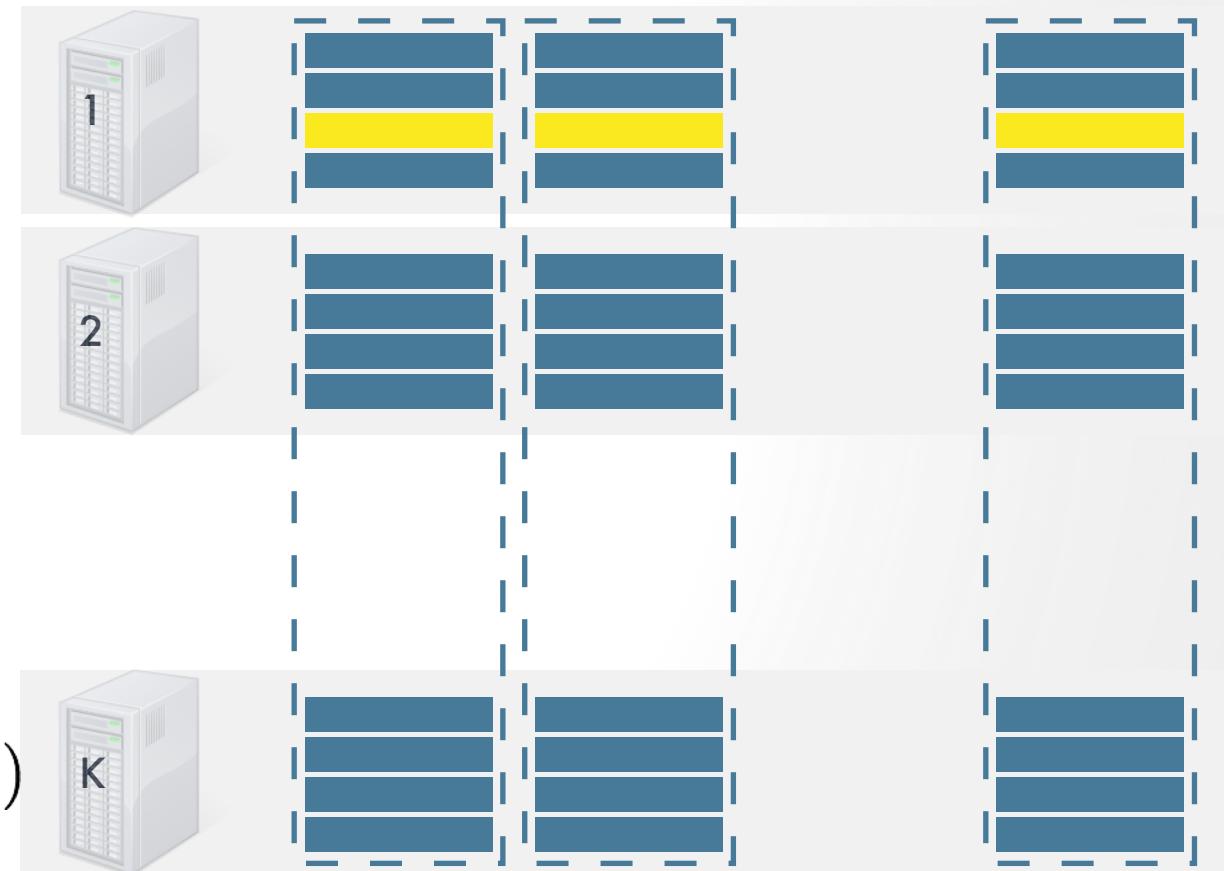


$\mathbf{math}(X_1)$

$\mathbf{math}(X_2)$

$\mathbf{math}(X_K)$

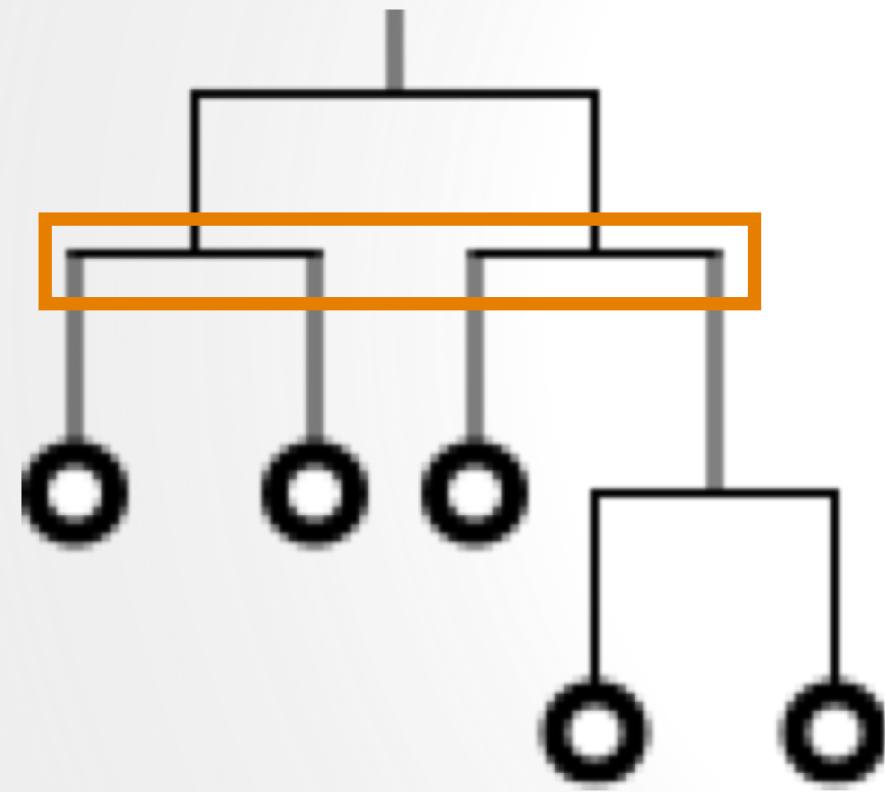
$\mathbf{X} = \{X_1, \dots, X_K\}$



# Tree Growth with Data Parallelism

$$\{X_i; t_i\} = f(\mathbf{math}(X_1), \dots, \mathbf{math}(X_K))$$

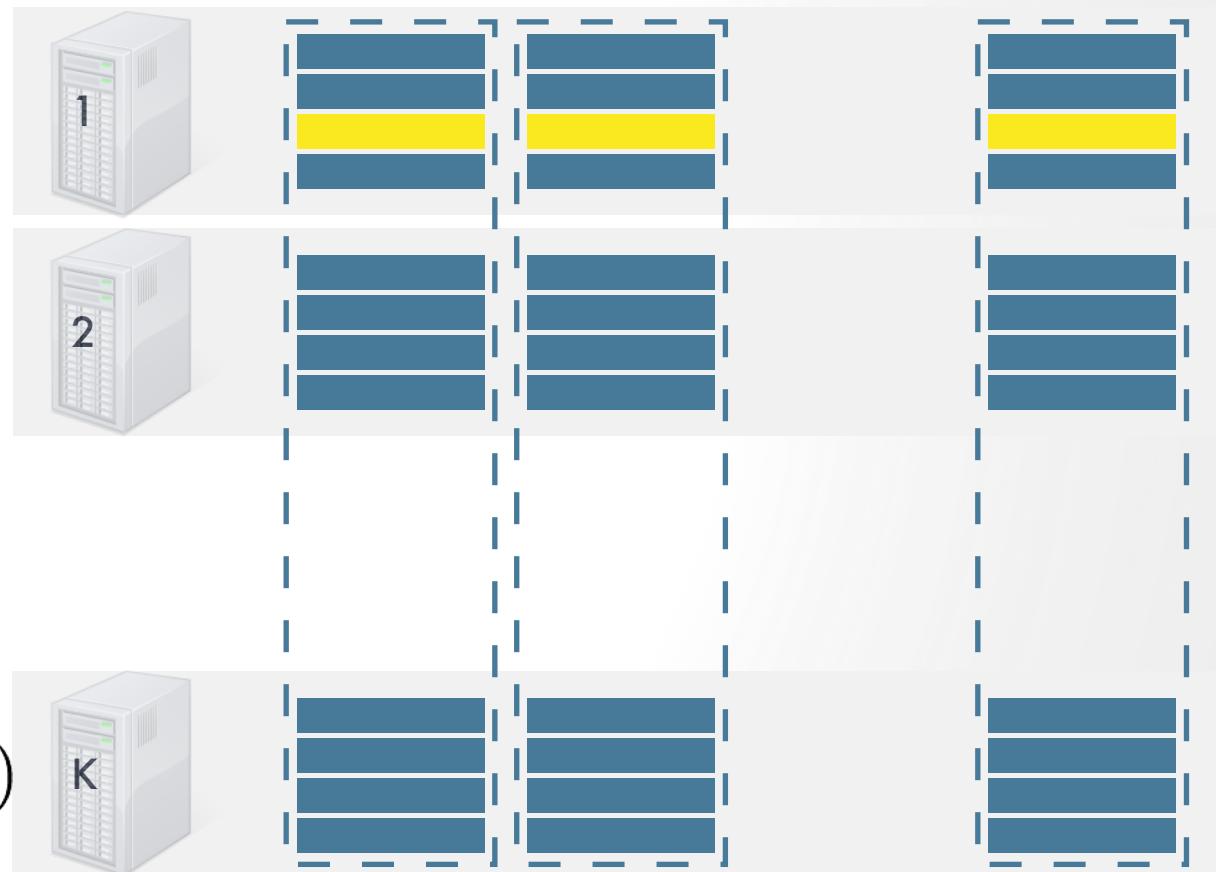
$$\mathbf{X} = \{X_1, \dots, X_K\}$$



$\mathbf{math}(X_1)$

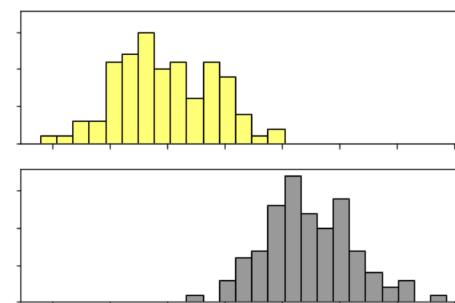
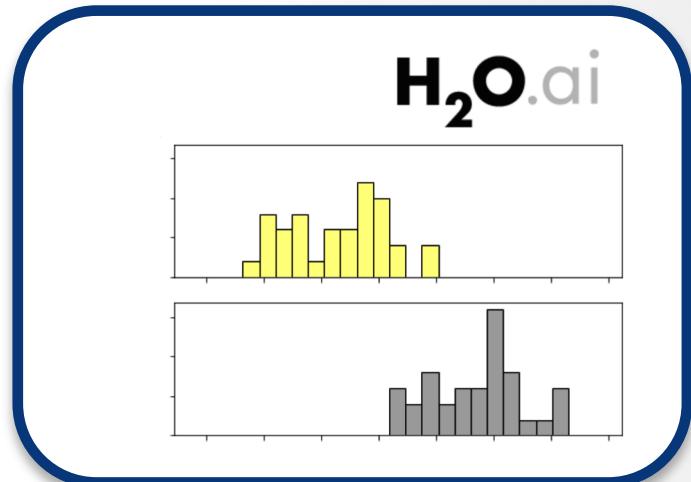
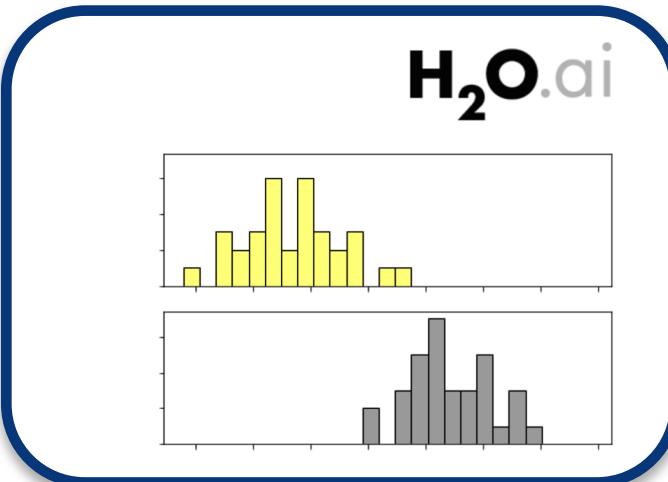
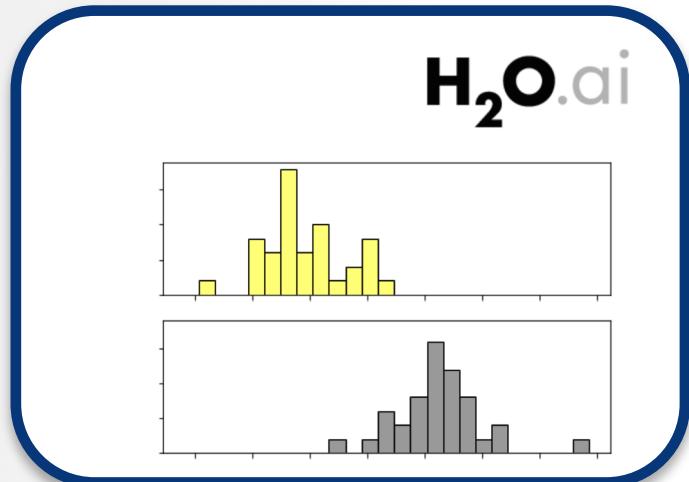
$\mathbf{math}(X_2)$

$\mathbf{math}(X_K)$

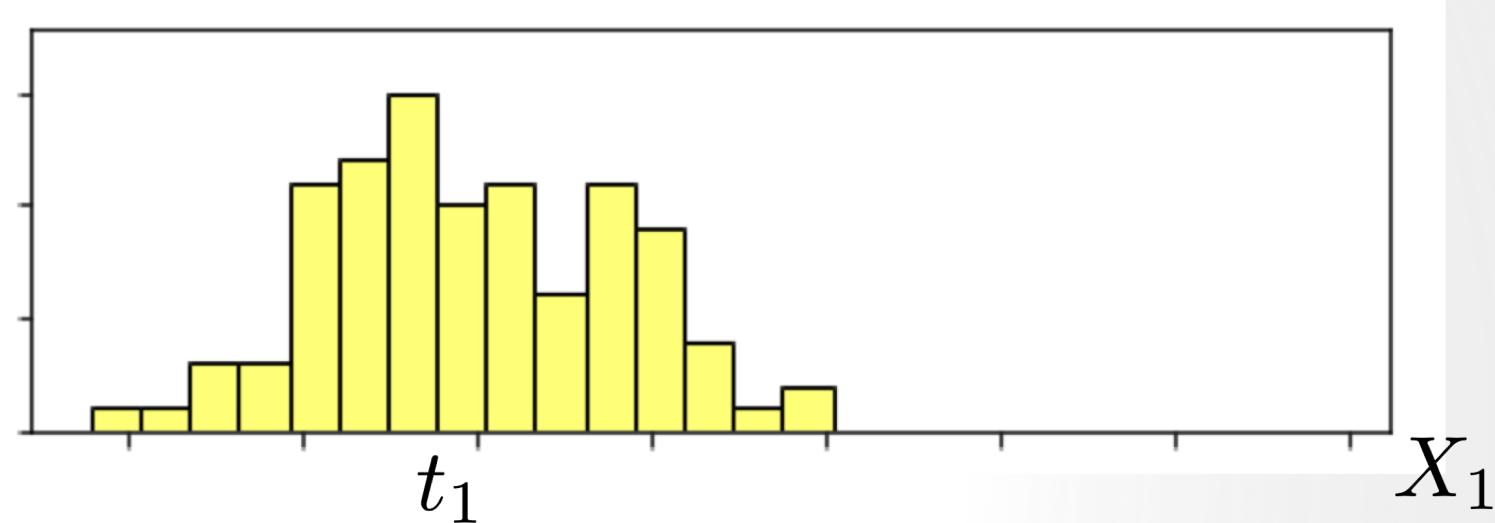
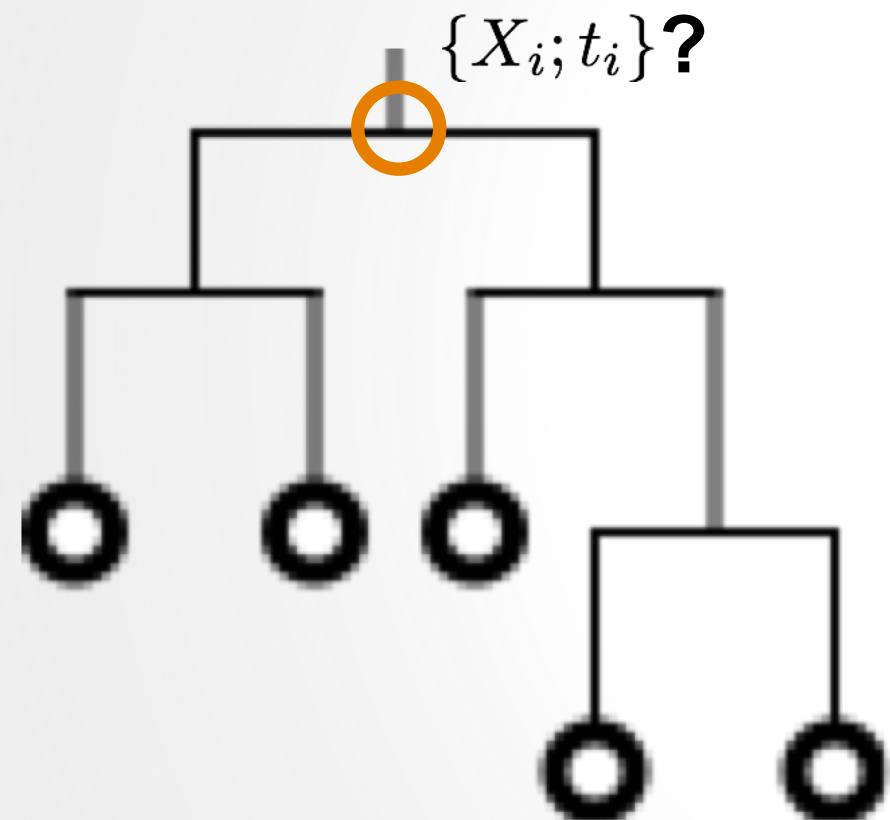


Full Data Parallelism for Each Level of Tree Growth!

# Splitting with Distributed Data



# Numerical Binning (Histogramming)



- **nbins**: number of bins in histogram
- **nbins\_top\_level**: number of bins to use at the top node, then halve at each ensuing level
- **histogram\_type**: method for binning {Uniform, Adaptive, Random, QuantilesGlobal, RoundRobin}

# Categorical Binning

- Lexigraphical ordering (i.e. alphabetical ordering)
- Example: {Red, Blue, Yellow, Orange, Purple, Green}
  - Lexigraphical order: {Blue, Green, Orange, Purple, Red, Yellow}
  - **nbin\_cats = 2**: {Blue, Green, Orange}, {Purple, Red, Yellow}
  - **nbin\_cats = 3**: {Blue, Green}, {Orange, Purple}, {Red, Yellow}
  - **nbin\_cats >= 6**: {Blue}, {Green}, {Orange}, {Purple}, {Red}, {Yellow}

# Binning in H2O Decision Trees

- Binning Numeric Features
  - Traditionally split points are chosen by sorting the each feature and inspecting an induced split.
  - For big data even when running parallel and distributed this can be computationally expensive so we approximate sorting with binning.
  - **More Bins, More Accurate** The number of bins can be specified by the user and it is the minimum number of bins required in a histogram built for each feature.
- Binning Categorical Features
  - **High Cardinality Features** slow model builds by inducing splits by each level.
  - Bin the levels in a categorical column according to “bins\_cat” parameter.
  - **More Bins, More Likely To Overfit** Increasing the number of bins can lead to splits on a single category, which can lead to overfitting.