# Decision Trees and other predictive models

## Mathias Lanner SAS Institute

# Agenda

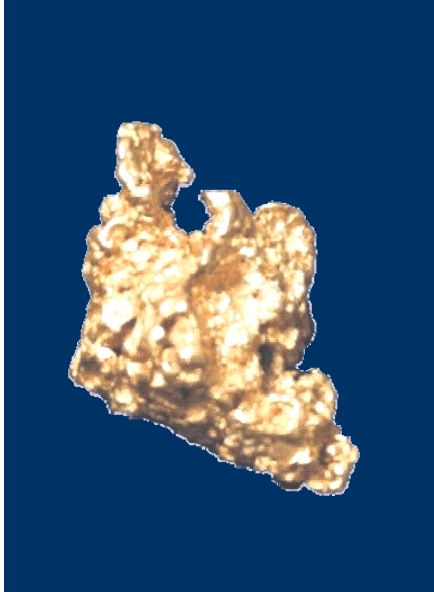| Introduction to Predictive Models |
|---|
| Decision Trees |
| Pruning |
| Regression |
| Neural Network |
| Model Assessment |
|  |

# Predictive Modeling

The *Essence* of Data Mining

*"Most of the big payoff [in data mining] has been in predictive modeling."*

– Herb Edelstein

# Predictive Modeling Applications

Database marketing

Financial risk management

Fraud detection

Process monitoring

Pattern detection

# Predictive Modeling Training Data

*Training Data*

| | | |
|---|---|---|
| *case 1:* | *inputs* | *target* |
| *case 2:* | *inputs* | *target* |
| *case 3:* | *inputs* | *target* |
| *case 4:* | *inputs* | *target* |
| *case 5:* | *inputs* | *target* |

**Numeric or categorical values**

# Predictive Modeling Score Data

### *Training Data*

| | | |
|---|---|---|
| *case 1:* | *inputs* | *target* |
| *case 2:* | *inputs* | *target* |
| *case 3:* | *inputs* | *target* |
| *case 4:* | *inputs* | *target* |
| *case 5:* | *inputs* | *target* |

### *Score Data*

| | | |
|---|---|---|
| *case 1:* | *inputs* | *?* |
| *case 2:* | *inputs* | *?* |
| *case 3:* | *inputs* | *?* |
| *case 4:* | *inputs* | *?* |
| *case 5:* | *inputs* | *?* |

**Only input values known**

# Predictions

### *Training Data*

| | | |
|---|---|---|
| *case 1:* | *inputs* | *target* |
| *case 2:* | *inputs* | *target* |
| *case 3:* | *inputs* | *target* |
| *case 4:* | *inputs* | *target* |
| *case 5:* | *inputs* | *target* |

### *Predictions*

### *Score Data*

| | | |
|---|---|---|
| *case 1:* | *inputs* | *?* |
| *case 2:* | *inputs* | *?* |
| *case 3:* | *inputs* | *?* |
| *case 4:* | *inputs* | *?* |
| *case 5:* | *inputs* | *?* |

# Predictions

### Training Data

| | | |
|---|---|---|
| case 1: | inputs | target |
| case 2: | inputs | target |
| case 3: | inputs | target |
| case 4: | inputs | target |
| case 5: | inputs | target |

### Predictions

prediction
prediction
prediction
prediction
prediction

### Score Data

| | | |
|---|---|---|
| case 1: | inputs | ? |
| case 2: | inputs | ? |
| case 3: | inputs | ? |
| case 4: | inputs | ? |
| case 5: | inputs | ? |

prediction
prediction
prediction
prediction
prediction

# Predictive Modeling Essentials

**Predict new cases**

**Select useful inputs**

**Optimize complexity**

# Predictive Modeling Essentials

**Predict new cases**

Select useful inputs

Optimize complexity

# Three Prediction Types

*Training Data*

| | | |
|---|---|---|
| *case 1:* | *inputs* | *target* |
| *case 2:* | *inputs* | *target* |
| *case 3:* | *inputs* | *target* |
| *case 4:* | *inputs* | *target* |
| *case 5:* | *inputs* | *target* |

*Predictions*

*prediction*
*prediction*
*prediction*
*prediction*
*prediction*

- **Decisions**

- **Rankings**

- **Estimates**

# Decision Predictions

**Training Data**

| | | |
|---|---|---|
| case 1: | inputs | target |
| case 2: | inputs | target |
| case 3: | inputs | target |
| case 4: | inputs | target |
| case 5: | inputs | target |

**Decisions**

primary
secondary
tertiary
primary
secondary

Trained model uses input measurements to make best decision for each case.

# Ranking Predictions

| *Training Data* | | *Rankings* | Trained model uses input measurements to optimally rank each case. |
|---|---|---|---|
| *case 1:* *inputs* | *target* | *720* | |
| *case 2:* *inputs* | *target* | *520* | |
| *case 3:* *inputs* | *target* | *620* | |
| *case 4:* *inputs* | *target* | *580* | |
| *case 5:* *inputs* | *target* | *470* | |

# Estimate Predictions

| *Training Data* | | | *Estimates* |
|---|---|---|---|
| *case 1:* | *inputs* | *target* | *0.65* |
| *case 2:* | *inputs* | *target* | *0.33* |
| *case 3:* | *inputs* | *target* | *0.54* |
| *case 4:* | *inputs* | *target* | *0.47* |
| *case 5:* | *inputs* | *target* | *0.28* |

**Trained model uses input measurements to optimally estimate target value.**

# Model Essentials – Predict Review

**Predict new cases**

**Decide, rank, estimate**

Select useful inputs

Optimize complexity

# Model Essentials – Select Review

 **Predict new cases**

 **Select useful inputs**

 **Optimize complexity**

# Curse of Dimensionality

1–D

2–D

3–D

# Input Selection

Redundancy                    Irrelevancy

# Model Essentials – Select Review

**Predict new cases**

**Decide, rank, estimate**

$x_1$ $x_2$ $x_3$ $x_4$ **Select useful inputs**

**Eradicate redundancies irrelevancies**

**Optimize complexity**

19

# Model Essentials – Optimize

Predict new cases

Select useful inputs

**Optimize complexity**

# Fool's Gold

My model fits the training data perfectly...

*I've struck it rich!*

# Model Complexity



Too flexible

Not flexible enough

# Data Splitting

# Training Data Role

## Training Data

| | | |
|---|---|---|
| case 1: | inputs | target |
| case 2: | inputs | target |
| case 3: | inputs | target |
| case 4: | inputs | target |
| case 5: | inputs | target |

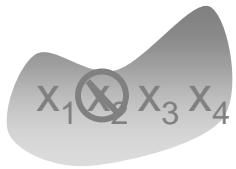*prediction*
*prediction*
*prediction*
*prediction*
*prediction*

**Training data gives sequence of predictive models with increasing complexity.**

## Validation Data

| | | |
|---|---|---|
| case 1: | inputs | target |
| case 2: | inputs | target |
| case 3: | inputs | target |
| case 4: | inputs | target |
| case 5: | inputs | target |

# Validation Data Role

*Training Data*

| | | |
|---|---|---|
| *case 1:* | *inputs* | *target* |
| *case 2:* | *inputs* | *target* |
| *case 3:* | *inputs* | *target* |
| *case 4:* | *inputs* | *target* |
| *case 5:* | *inputs* | *target* |

*prediction*
*prediction*
*prediction*
*prediction*
*prediction*

*Validation Data*

| | | |
|---|---|---|
| *case 1:* | *inputs* | *target* |
| *case 2:* | *inputs* | *target* |
| *case 3:* | *inputs* | *target* |
| *case 4:* | *inputs* | *target* |
| *case 5:* | *inputs* | *target* |

**Validation data helps select best model from sequence**

# Validation Data Role

*Training Data*

*case 1:  inputs        target*
*case 2:  inputs        target*
*case 3:  inputs        target*
*case 4:  inputs        target*
*case 5:  inputs        target*

*prediction*
*prediction*
*prediction*
*prediction*
*prediction*

## Validation Data

*case 1:  inputs        target*
*case 2:  inputs        target*
*case 3:  inputs        target*
*case 4:  inputs        target*
*case 5:  inputs        target*

*prediction*
*prediction*
*prediction*
*prediction*
*prediction*

**Validation data helps select best model from Sequence.**

# Model Essentials – Optimize

**Predict new cases**

Decide, rank, estimate

**Select useful inputs**

$x_1$ $x_2$ $x_3$ $x_4$

Eradicate redundancies irrelevancies

**Optimize complexity**

**Tune models with validation data**

# Agenda

| |
|---|
| **Introduction to Predictive Models** |
| **DECISION TREES** |
| **Pruning** |
| **Regression** |
| **Neural Networks** |
| **Model Assessment** |
| |

28

# Predictive Modeling Tools

| Primary | Decision Tree | Regression | Neural Network |
|---------|---------------|------------|----------------|

| Specialty | Dmine Regression | MBR | AutoNeural |
|-----------|------------------|-----|------------|
| | Rule Induction | DMNeural | |

| Multiple Model | Ensemble | Two Stage |
|----------------|----------|-----------|

# Predictive Modeling Tools

**Primary**

Decision Tree    Regression    Neural Network

**Specialty**

Dmine Regression    MBR    AutoNeural

Rule Induction    DMNeural

**Multiple Model**

Ensemble    Two Stage

# Predictive Modeling Tools

**Primary**

Decision Tree    Regression    Neural Network

**Specialty**

Dmine Regression    MBR    AutoNeural

Rule Induction    DMNeural

**Multiple Model**

Ensemble    Two Stage

# Model Essentials – Decision Trees

**Predict new cases** — **Prediction rules**

**Select useful inputs** — **Split search**

$x_1 \ x_2 \ x_3 \ x_4$

**Optimize complexity** — **Pruning**

# Simple Prediction Illustration

**Analysis goal:**

**Predict the color of a dot based on its location in a scatter plot.**

# Model Essentials – Decision Trees

**Predict new cases**

**Prediction rules**

Select useful inputs

Split search

Optimize complexity

Pruning

# Decision Tree Prediction Rules

# Decision Tree Prediction Rules

# Decision Tree Prediction Rules

# Model Essentials – Decision Trees

**Predict new cases**            **Prediction rules**

**Select useful inputs**         **Split search**

**Optimize complexity**          **Pruning**

# Decision Tree Split Search

*left*     *right*



**Calculate the *logworth* of every partition on input $x_1$.**

# Decision Tree Split Search



*left*  *right*

● 53% | 42%

max logworth($x_1$)
0.95

● 47% | 58%

**Select the partition with maximum *logworth*.**

40

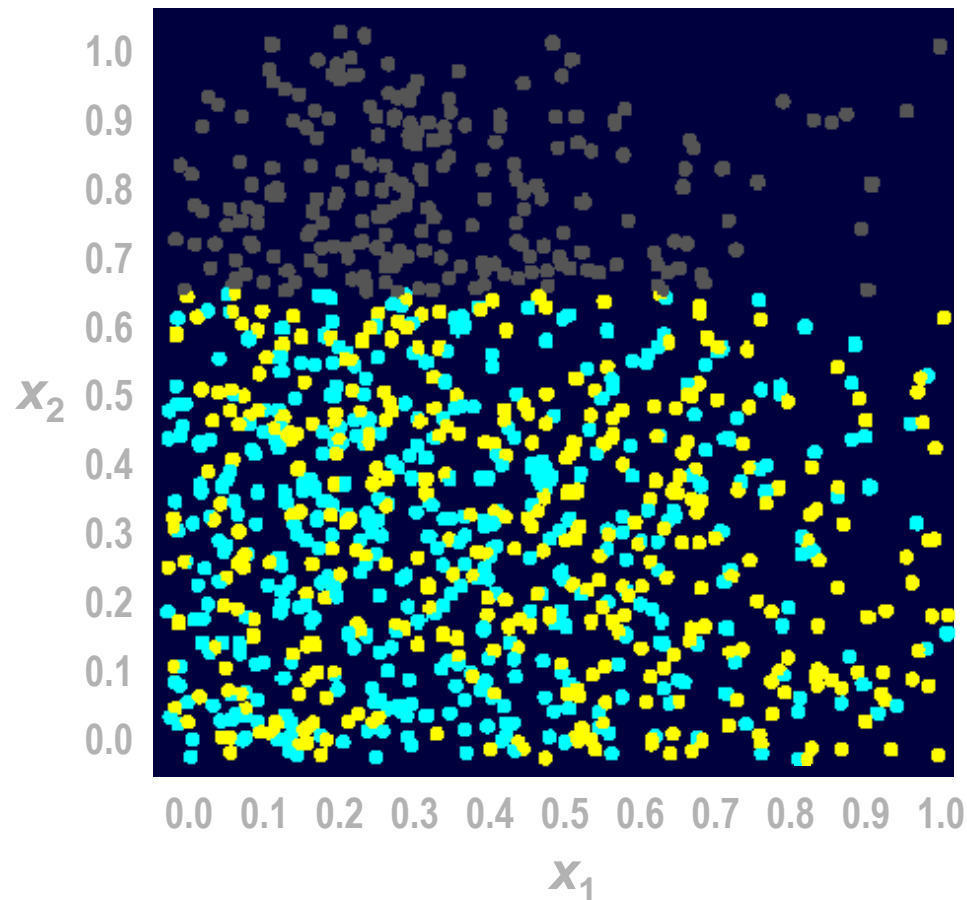# Decision Tree Split Search

# Decision Tree Split Search

# Decision Tree Split Search



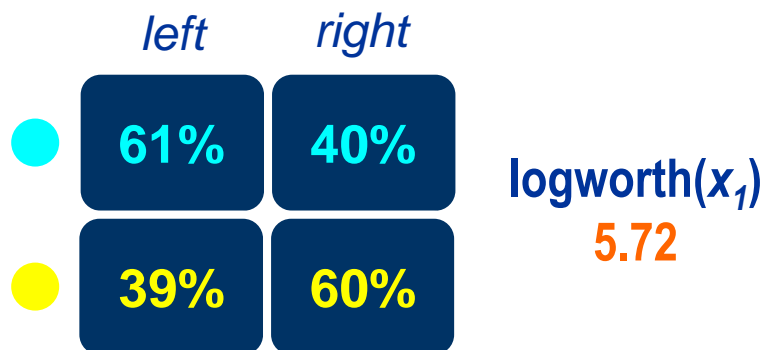Create partition rule from best partition across all inputs.

# Decision Tree Split Search



**Repeat process in each subset.**
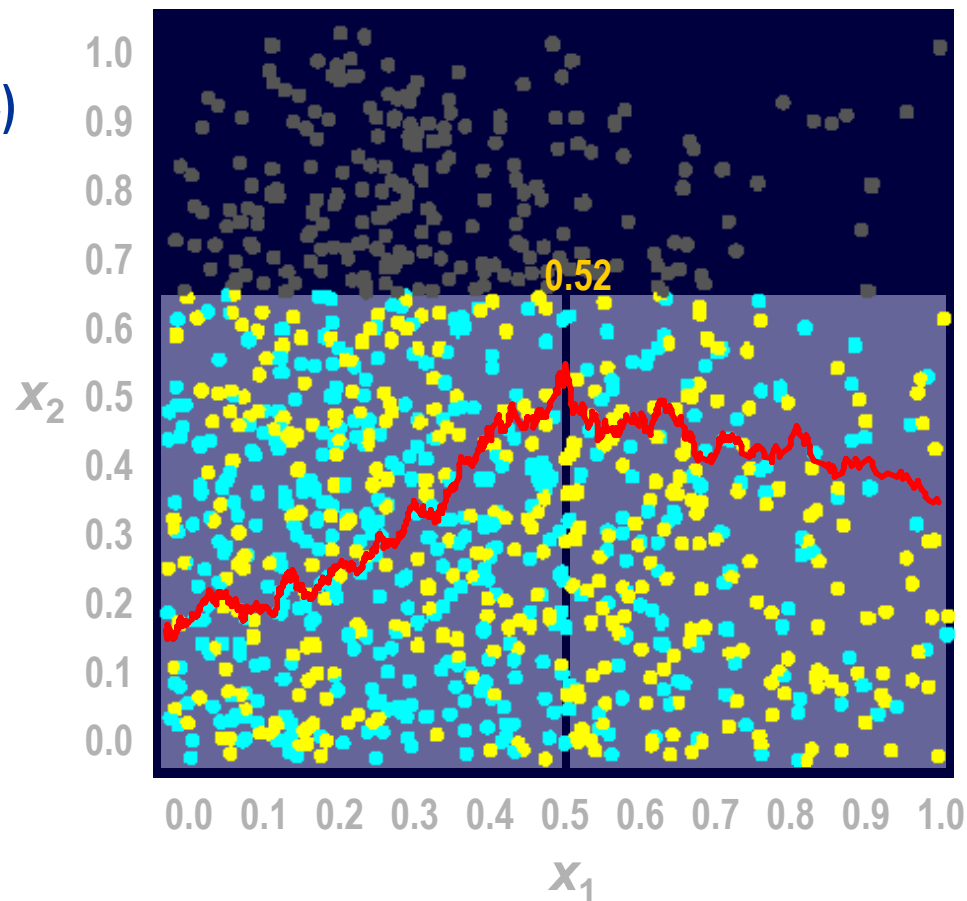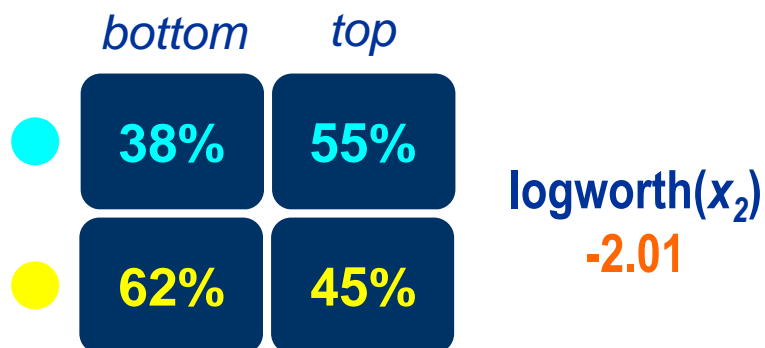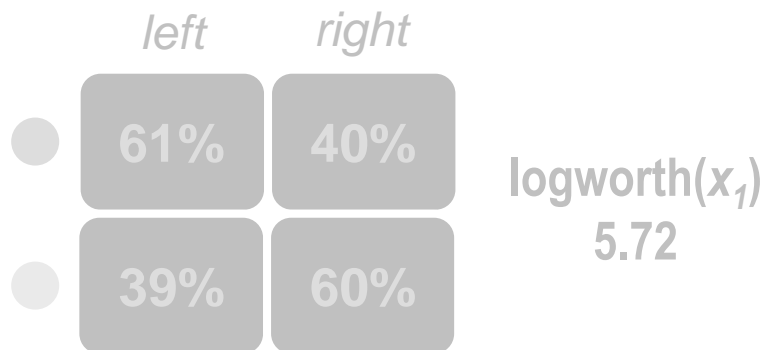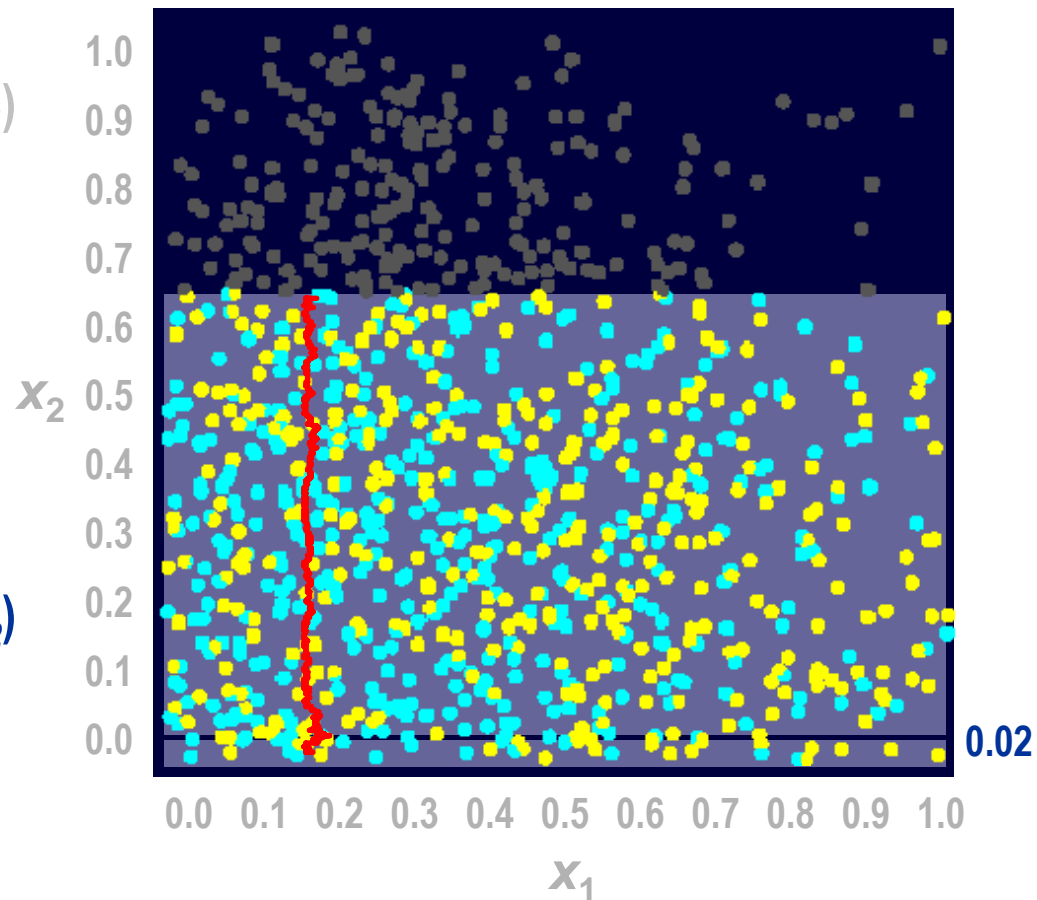
# Decision Tree Split Search



left    right

61%    40%

logworth($x_1$)
5.72

39%    60%

**Select the partition with maximum *logworth* on input $x_1$.**

# Decision Tree Split Search
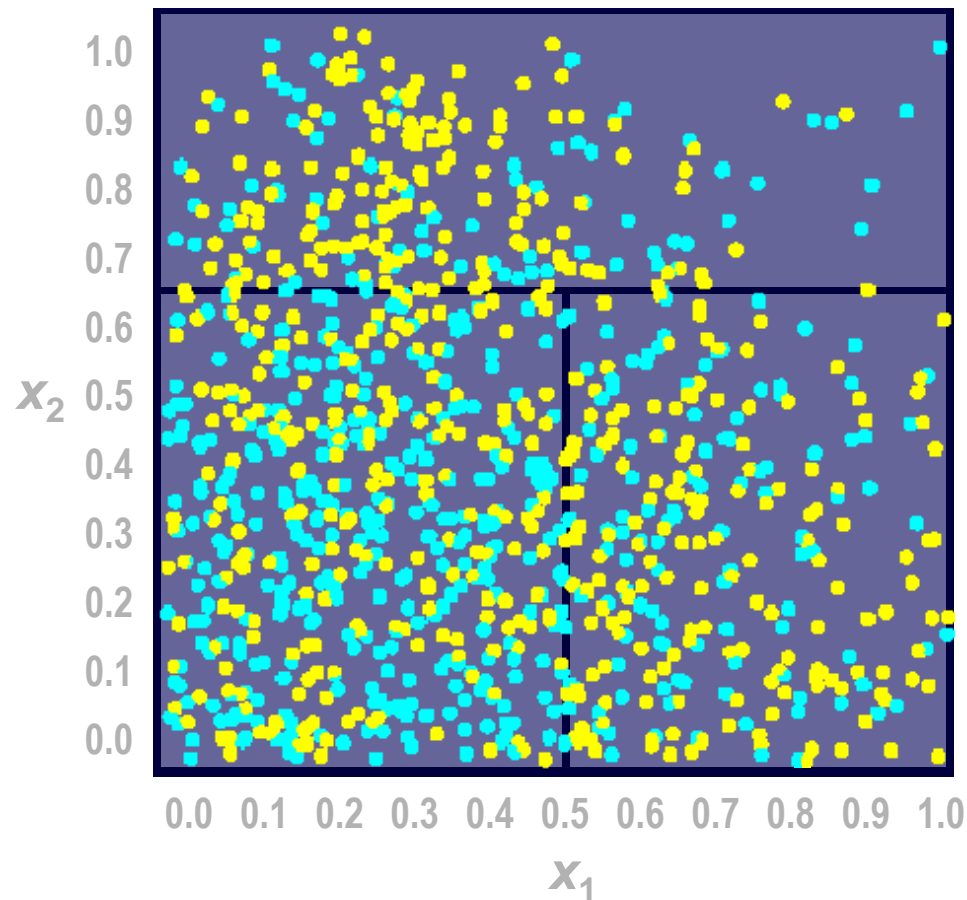
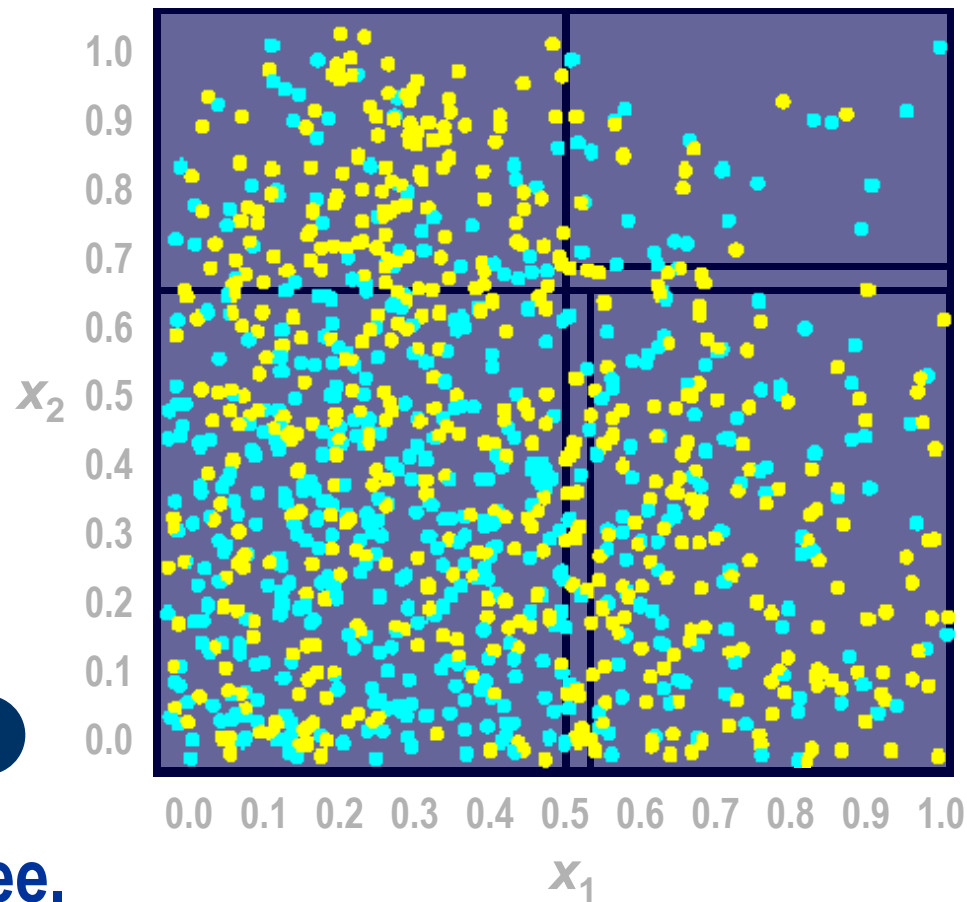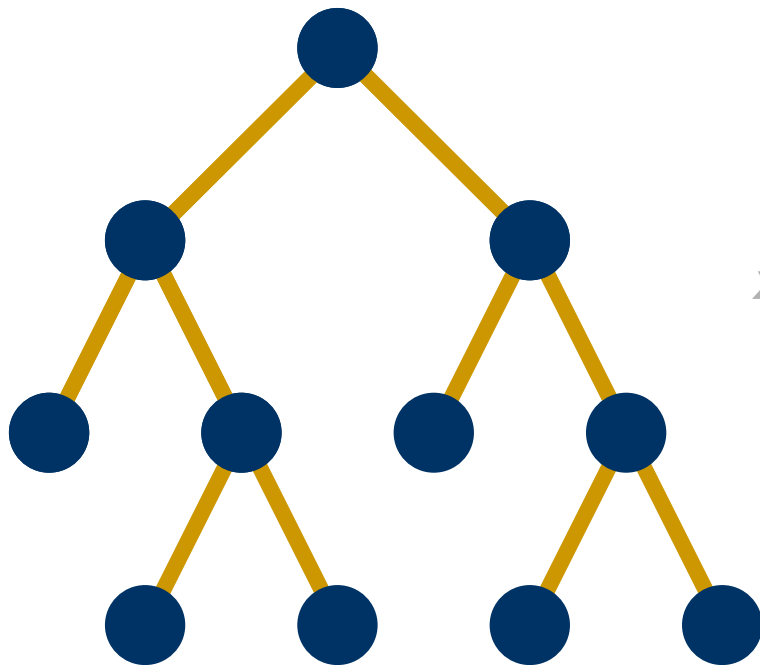# Decision Tree Split Search



**Create second partition rule.**

# Decision Tree Split Search



**Repeat to form maximal tree.**

# Demo Decision Tree