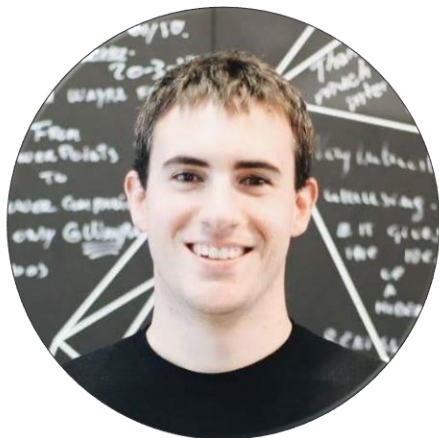# DISTRIBUTED LOGISTIC MODEL TREES

## 16 NOV 2016 @ APACHE BIG DATA EUROPE

Distributed Logistic Model Trees, Stratio Intelligence

**Mateo Álvarez** and **Antonio Soriano**

@StratioBD

# PROFILE

## SKILLS



## MATEO ÁLVAREZ

Aerospace Engineer, MSc in
Propulsion Systems (UPM), Master
in Data Science (URJC).
Working as data scientist and Big
Data developer at Stratio Big Data in
the data science department

**in** mateo-alvarez

# PROFILE

## ANTONIO SORIANO

Ph.D. in Telecommunications, MSc in Electronic Systems Engineering and Telecommunication Technologies, Systems and Networks (UPV), and MSc "Big Data Expert" (UTAD).
Working as data scientist and Big Data developer at at Stratio Big Data in the data science department

@Phd_A_Soriano

## SKILLS

R

Scala

python

Spark

Java

# INDEX

# INTERPRETABLE ALGORITHMS

- Why use interpretable algorithms instead of "black boxes"
- Logistic Regression
- Decision Trees
- Variance-Bias tradeoff
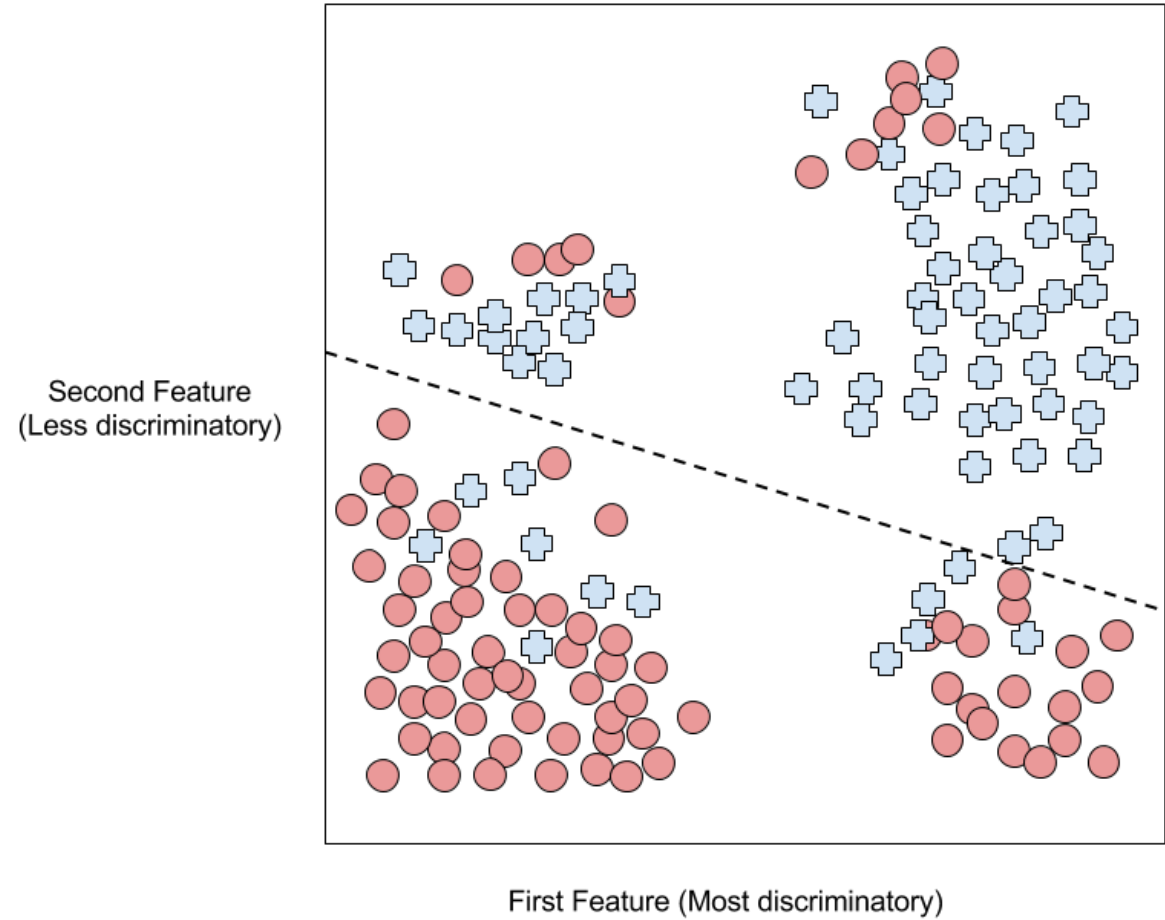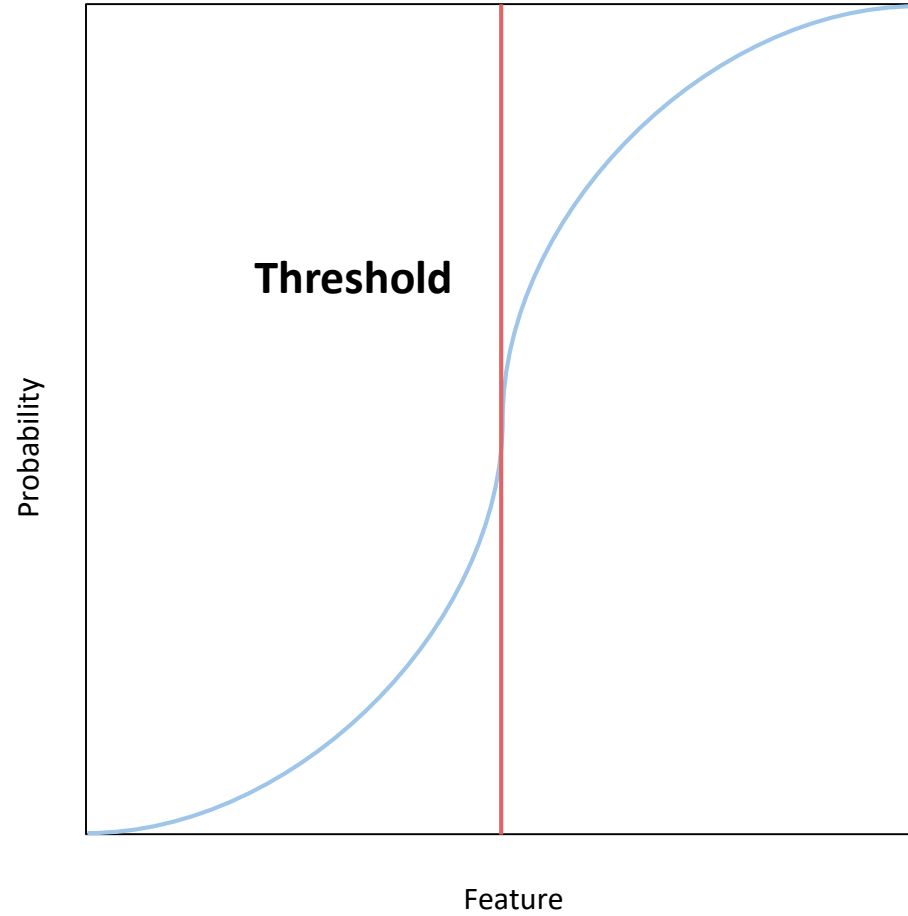
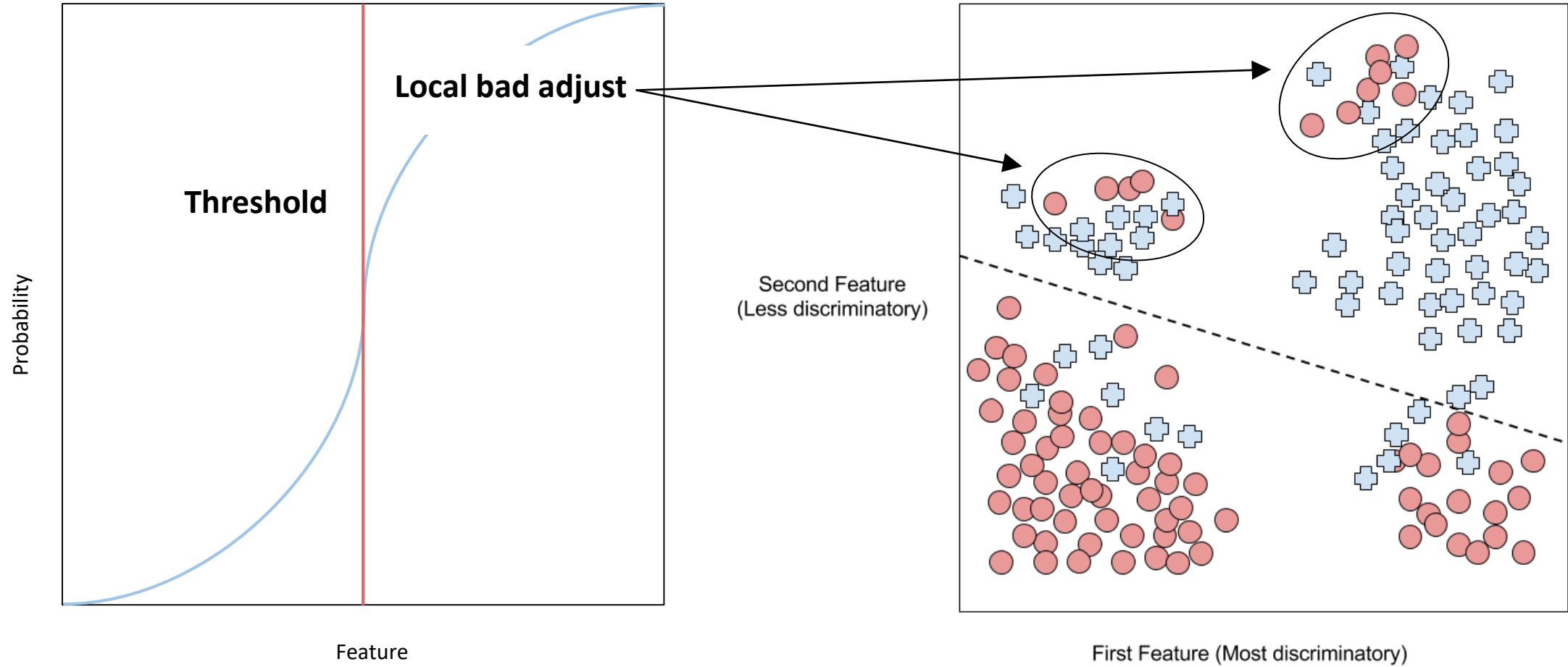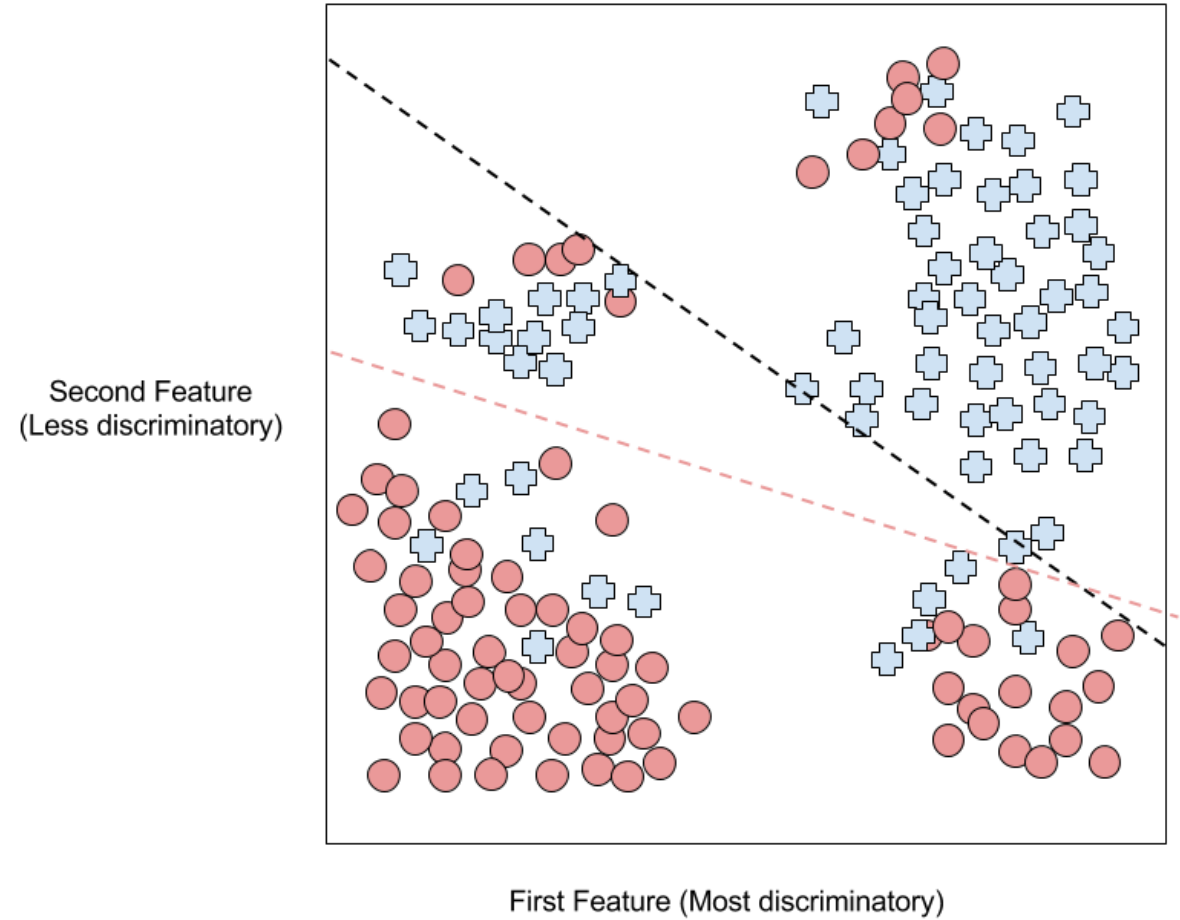@StratioBD
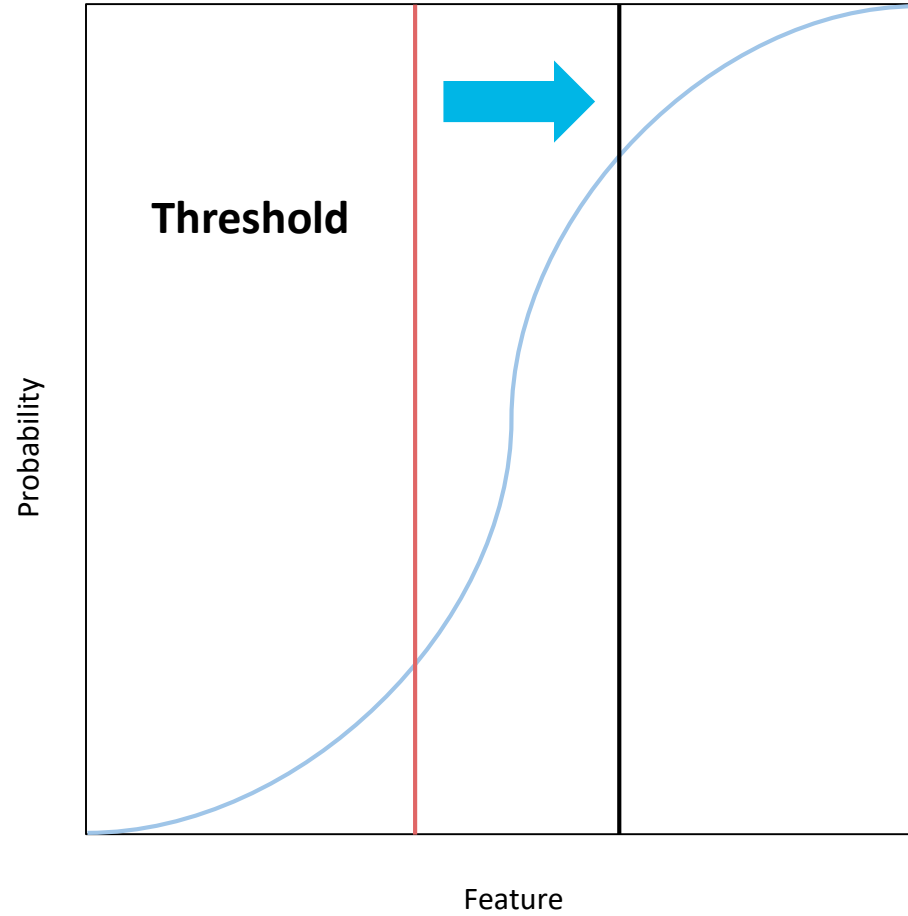
Accuracy **VS** Explainability

Medical Studies  Power management  Financial environment  Criminal activity
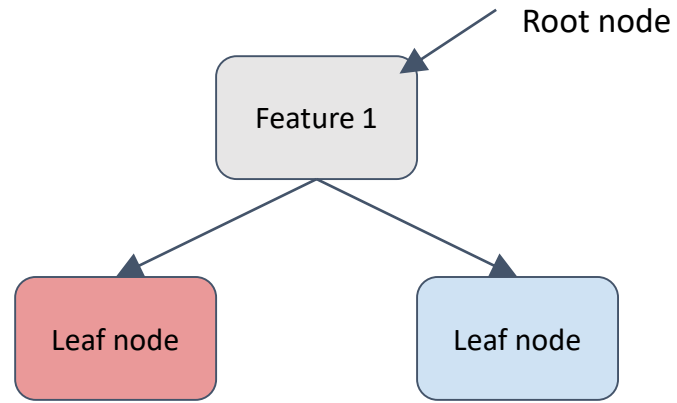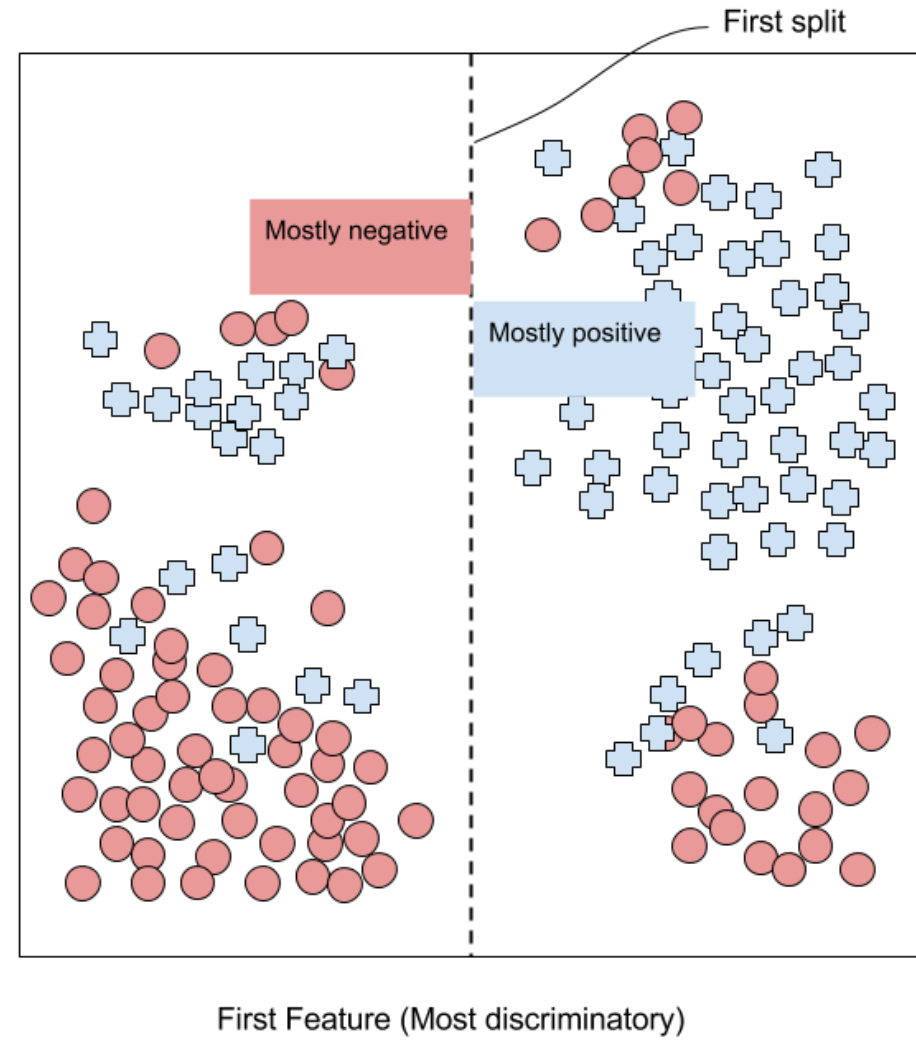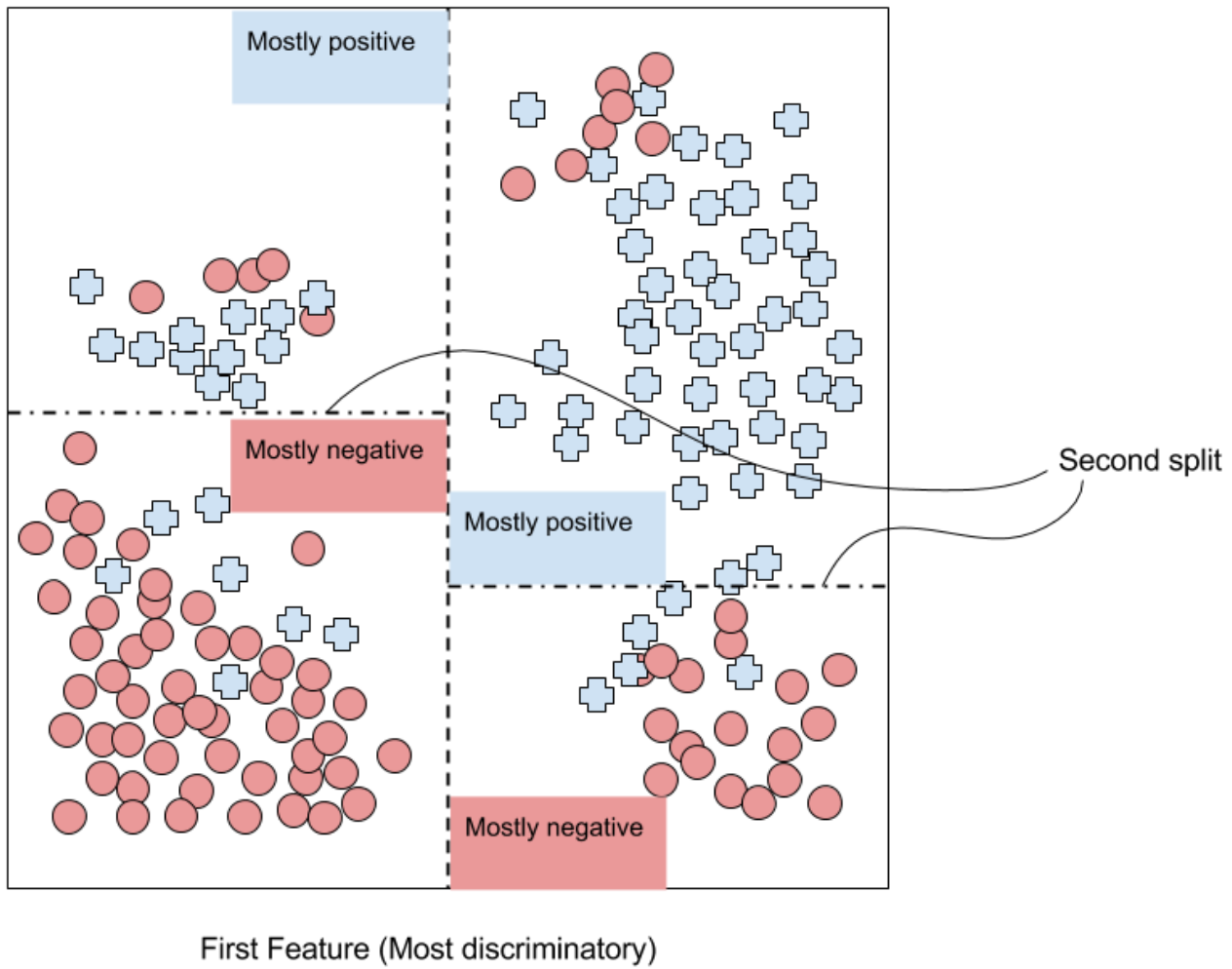
# LOGISTIC REGRESSION

# LOGISTIC REGRESSION

# DECISION TREES

Root node

Feature 1

Leaf node

Leaf node

First split

Mostly negative

Mostly positive

Second Feature
(Less discriminatory)

First Feature (Most discriminatory)

# DECISION TREES



Root node

Feature 1

Feature 2

Feature 2

Leaf node

Leaf node

Leaf node

Leaf node
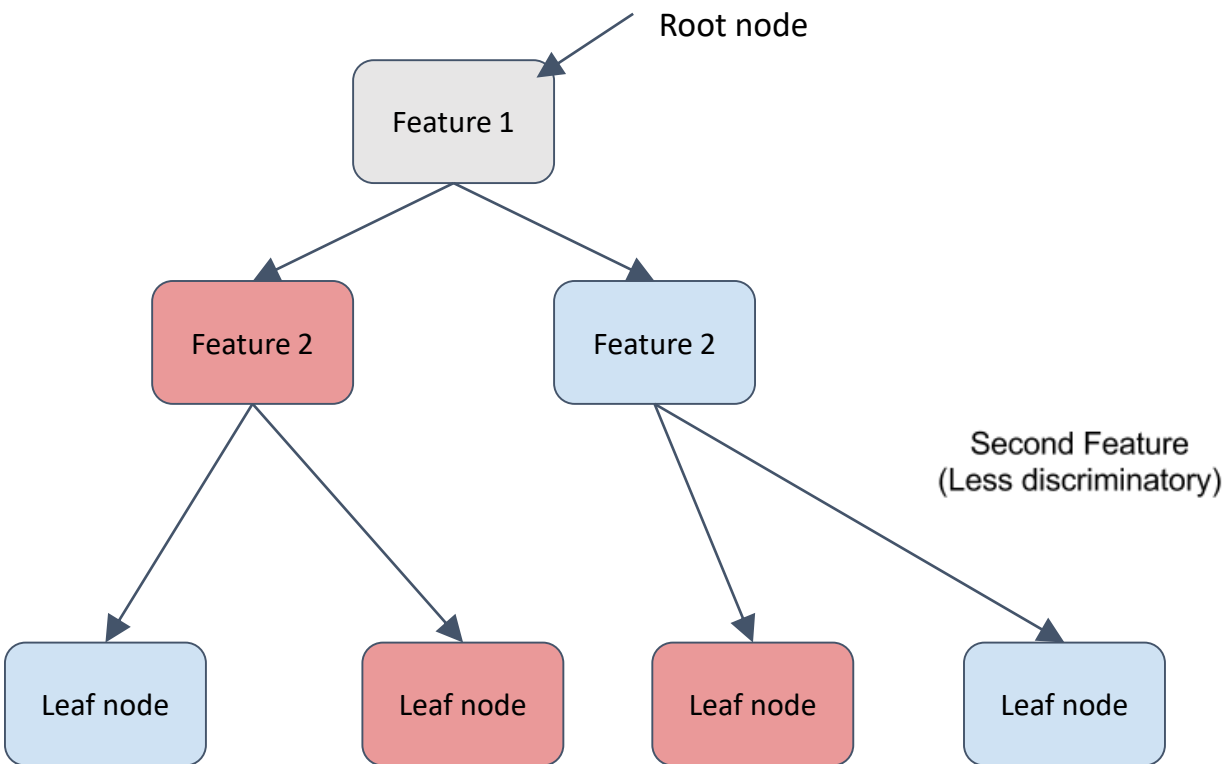
Second Feature
(Less discriminatory)

Mostly positive

Mostly negative

Mostly positive

Mostly negative

Second split

First Feature (Most discriminatory)

# DECISION TREES

**Test error**

**Training error**

Mean Error

Model complexity

**Total error**

Underfitting Overfitting

**Variance**

**Bias²**

Error

Model complexity

# BIAS-VARIANCE TRADE-OFF

$$\mathrm{E}\left[\left(y - \hat{f}(x)\right)^2\right] = \mathrm{Bias}\left[\hat{f}(x)\right]^2 + \mathrm{Var}\left[\hat{f}(x)\right] + \sigma^2$$

Where:

$$\mathrm{Bias}\left[\hat{f}(x)\right] = \mathrm{E}\left[\hat{f}(x) - f(x)\right]$$

and

$$\mathrm{Var}\left[\hat{f}(x)\right] = \mathrm{E}[\hat{f}(x)^2] - \mathrm{E}[\hat{f}(x)]^2$$

Bias

Variance

# BIAS-VARIANCE TRADE-OFF

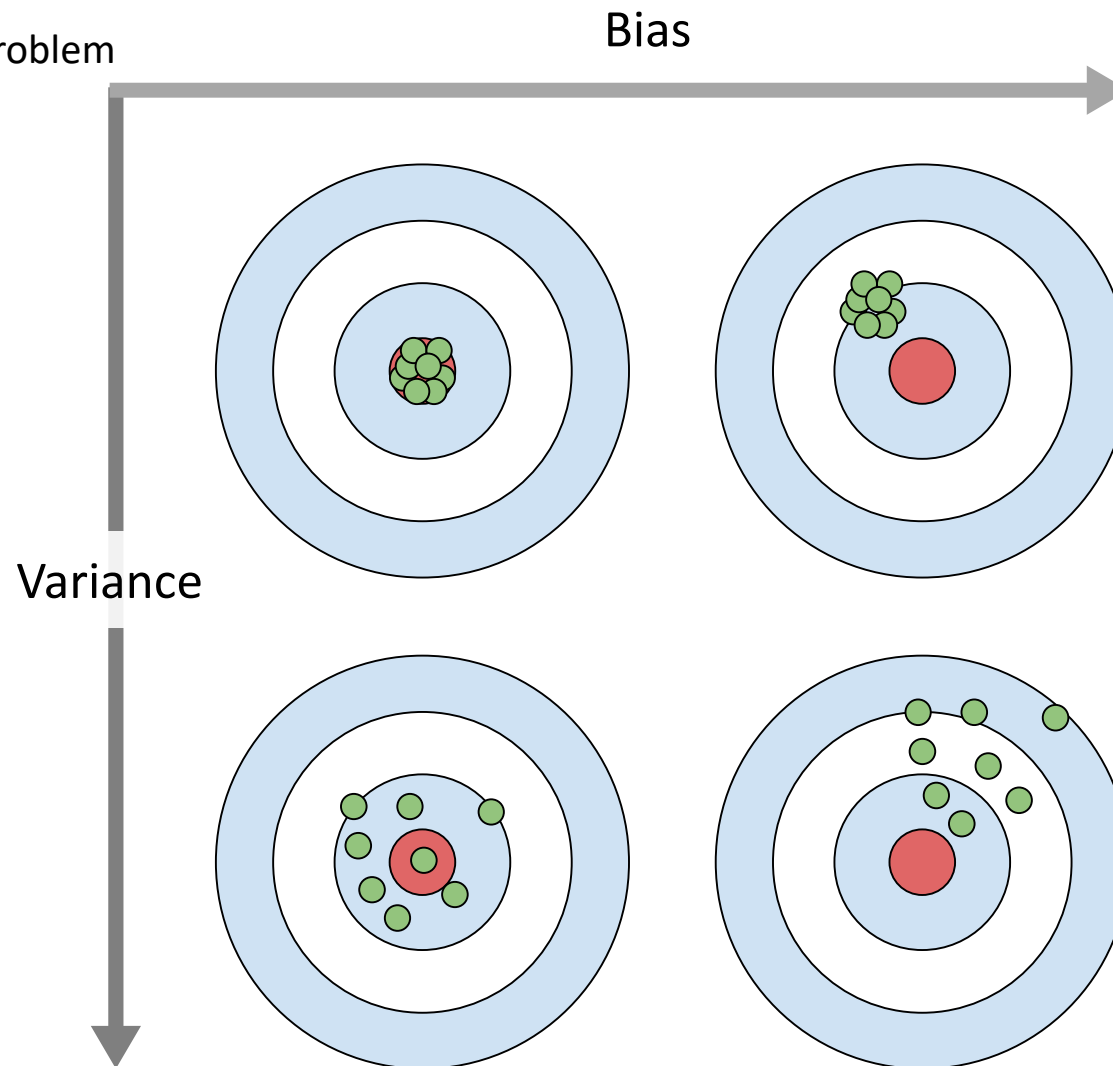Missing important variables for the problem to make the predictions

$$\mathrm{E}\left[(y - \hat{f}(x))^2\right] = \mathrm{Bias}\left[\hat{f}(x)\right]^2 + \mathrm{Var}\left[\hat{f}(x)\right] + \sigma^2$$

Where:

$$\mathrm{Bias}\left[\hat{f}(x)\right] = \mathrm{E}\left[\hat{f}(x) - f(x)\right]$$

and

$$\mathrm{Var}\left[\hat{f}(x)\right] = \mathrm{E}[\hat{f}(x)^2] - \mathrm{E}[\hat{f}(x)]^2$$

Bias

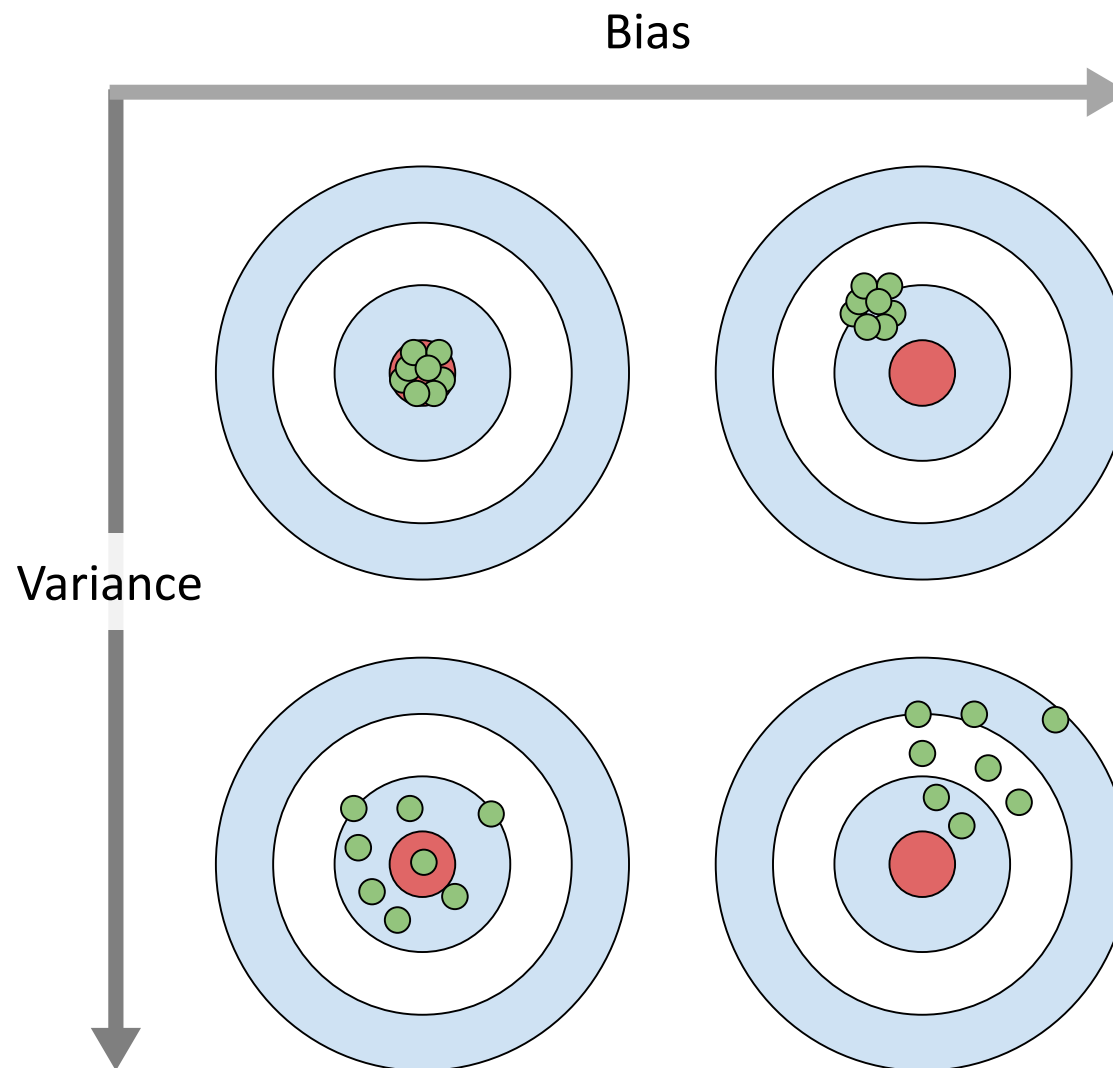Variance

Overfitting to the sample/training data

$$\mathrm{E}\left[\left(y - \hat{f}(x)\right)^2\right] = \mathrm{Bias}\left[\hat{f}(x)\right]^2 + \mathrm{Var}\left[\hat{f}(x)\right] + \sigma^2$$

Where:

$$\mathrm{Bias}\left[\hat{f}(x)\right] = \mathrm{E}\left[\hat{f}(x) - f(x)\right]$$

and

$$\mathrm{Var}\left[\hat{f}(x)\right] = \mathrm{E}[\hat{f}(x)^2] - \mathrm{E}[\hat{f}(x)]^2$$



Bias

Variance

Irreducible error on prediction

$$\mathrm{E}\left[\left(y - \hat{f}(x)\right)^2\right] = \mathrm{Bias}\left[\hat{f}(x)\right]^2 + \mathrm{Var}\left[\hat{f}(x)\right] + \sigma^2$$

Where:

$$\mathrm{Bias}\left[\hat{f}(x)\right] = \mathrm{E}\left[\hat{f}(x) - f(x)\right]$$

and

$$\mathrm{Var}\left[\hat{f}(x)\right] = \mathrm{E}[\hat{f}(x)^2] - \mathrm{E}[\hat{f}(x)]^2$$
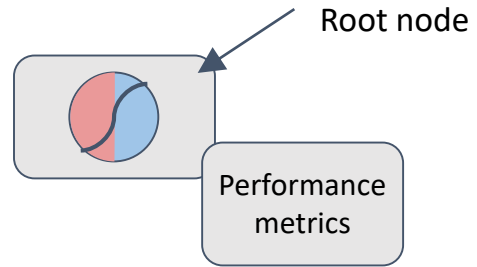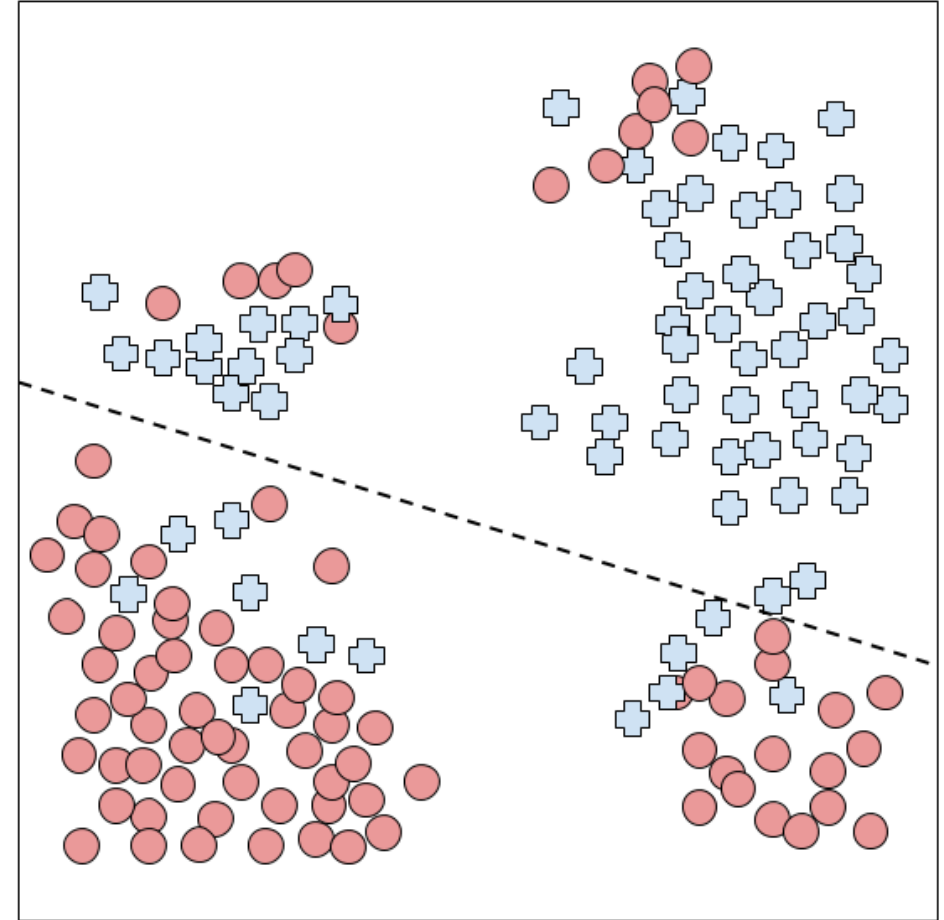
Bias

Variance

# DISTRIBUTED LOGISTIC MODEL TREES

- Logistic Model Trees

- Distributed implementation

- Cost function & configuration parametres

- Demo

@StratioBD

# LOGISTIC MODEL TREES

# LOGISTIC MODEL TREES

# LOGISTIC MODEL TREES



Root node

Feature 1

Performance metrics

Performance metrics

Performance metrics

Performance metrics

Performance metrics

Performance metrics

Second Feature (Less discriminatory)

Logistic Regression splits

First Feature (Most discriminatory)
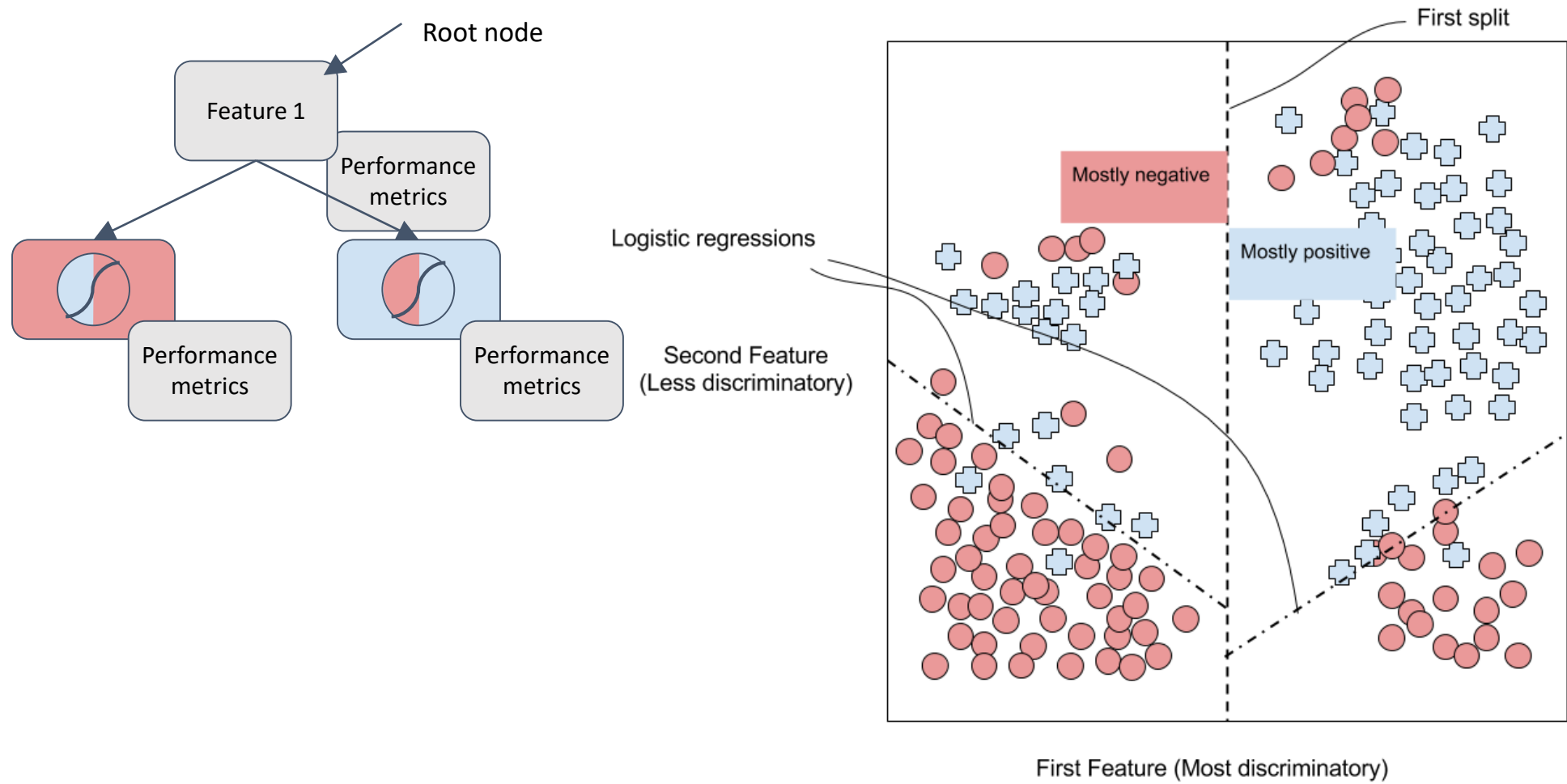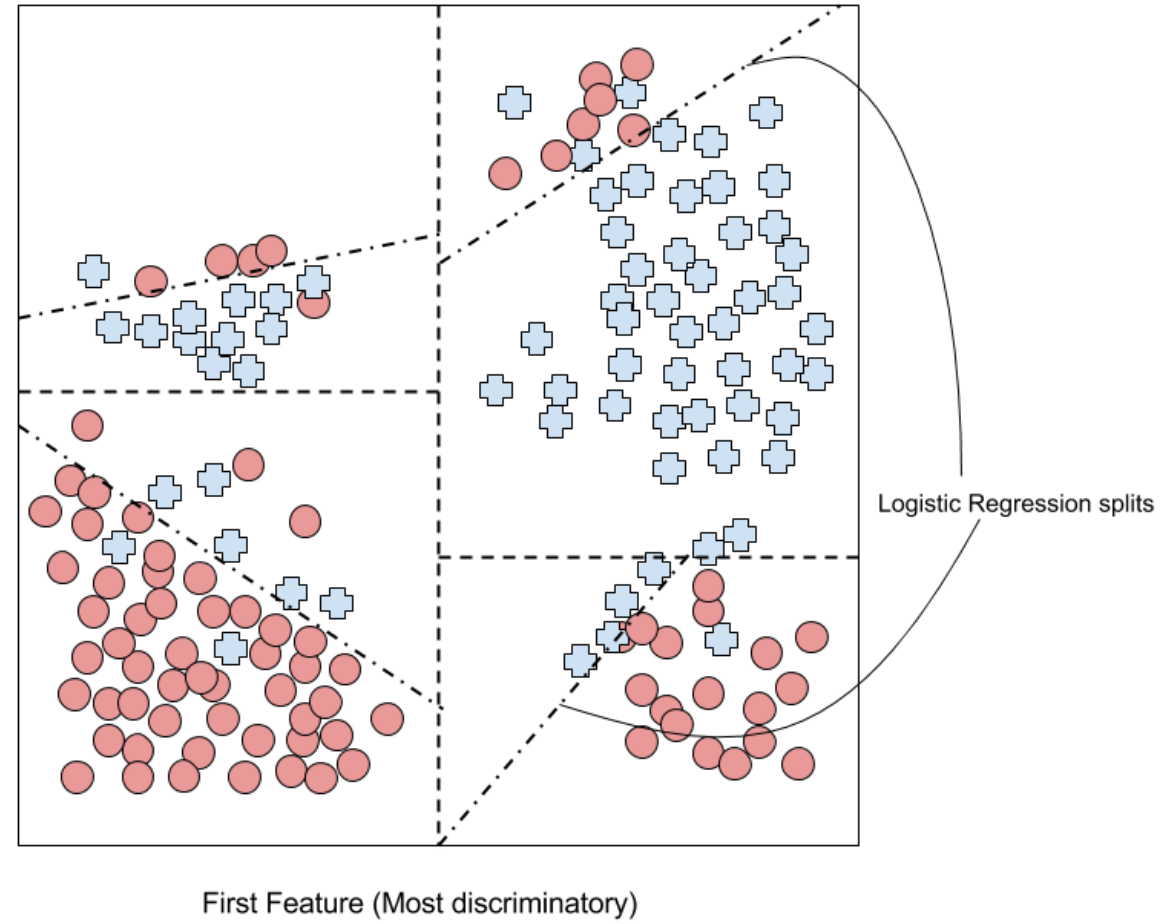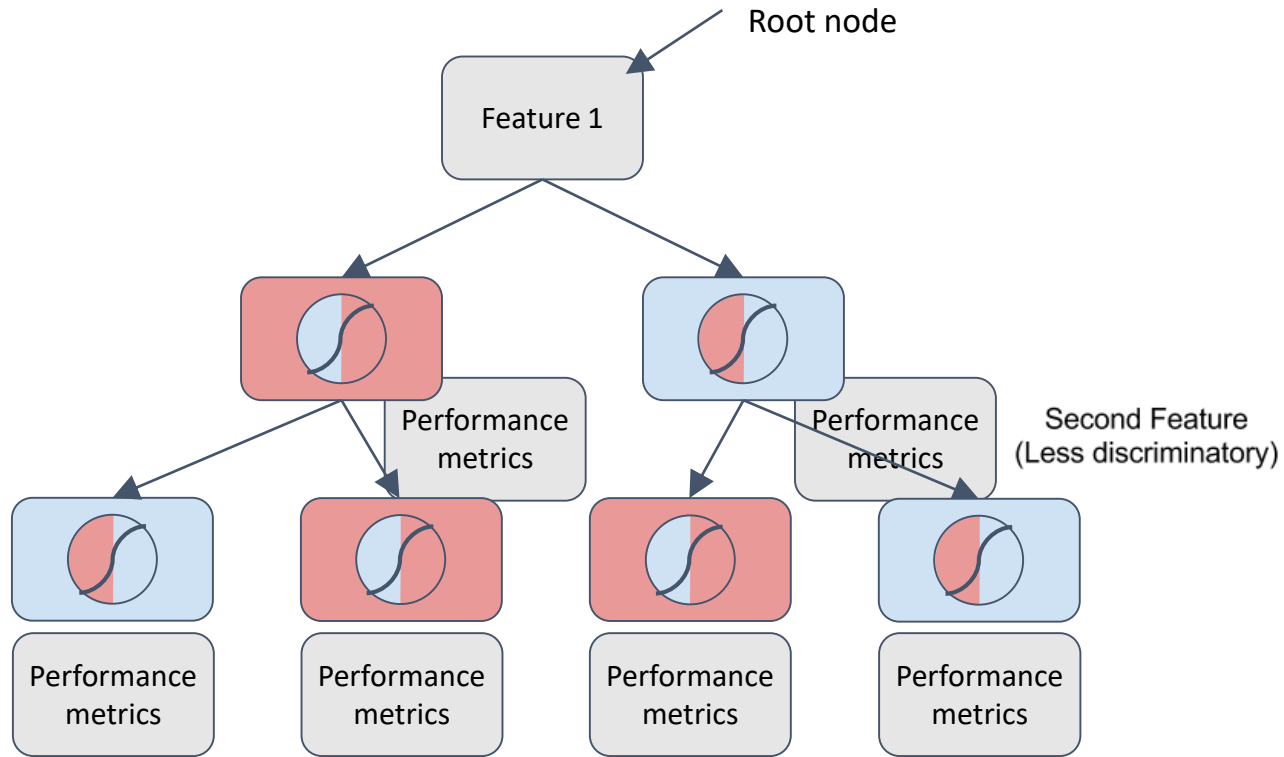
# LOGISTIC MODEL TREES

Root node

Feature 1

Feature 2

Second Feature
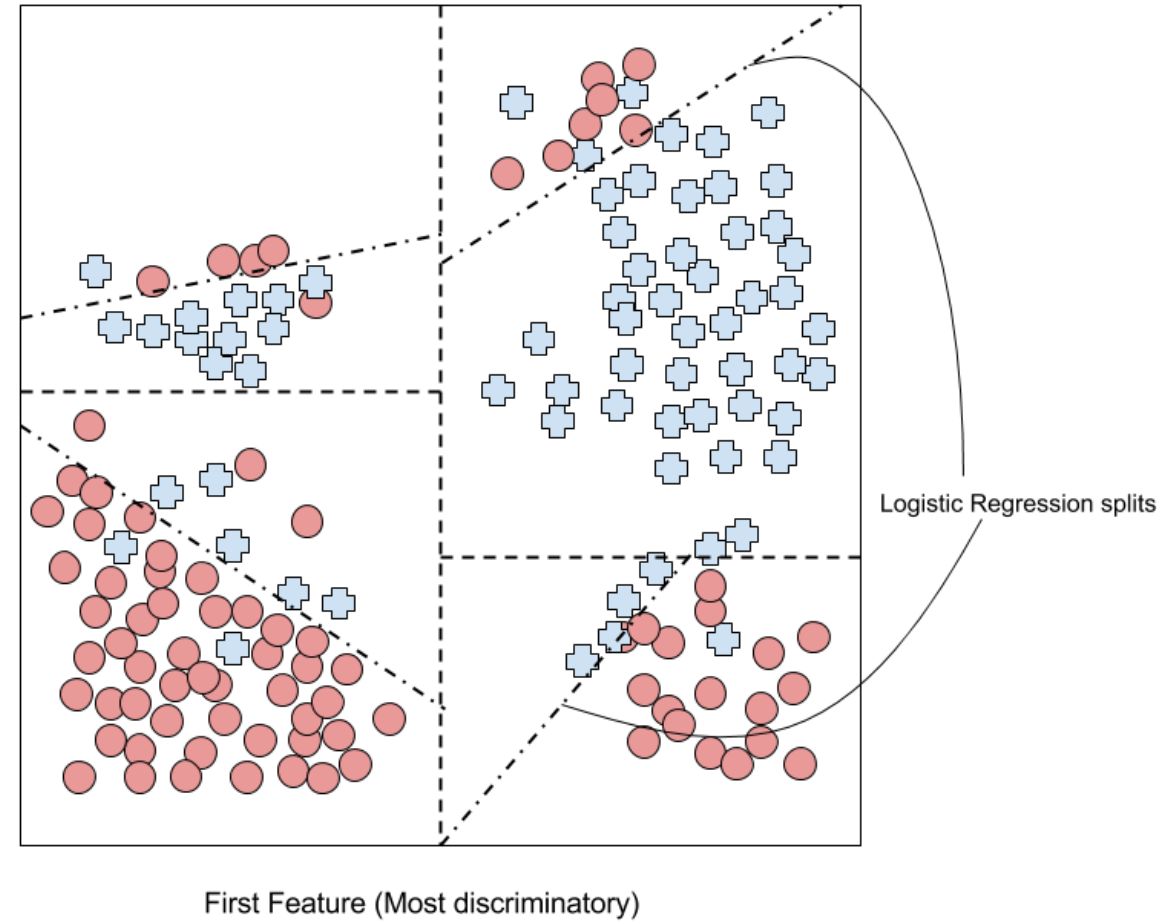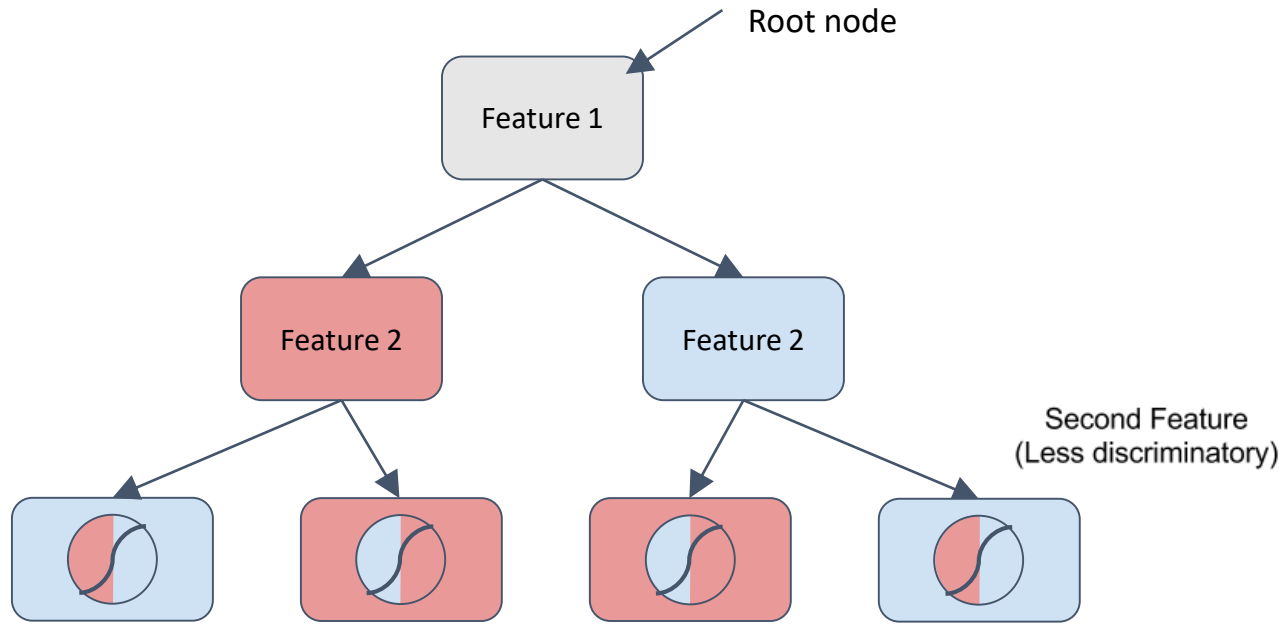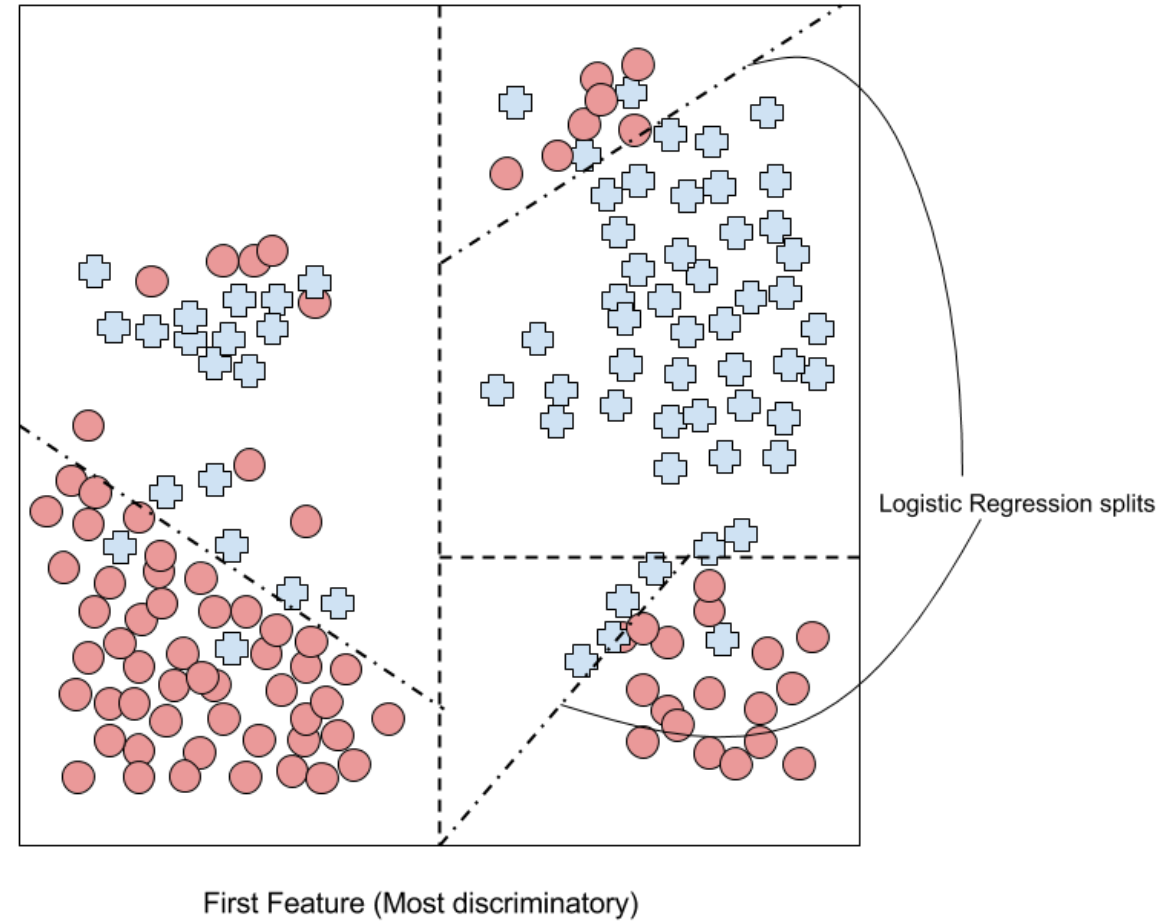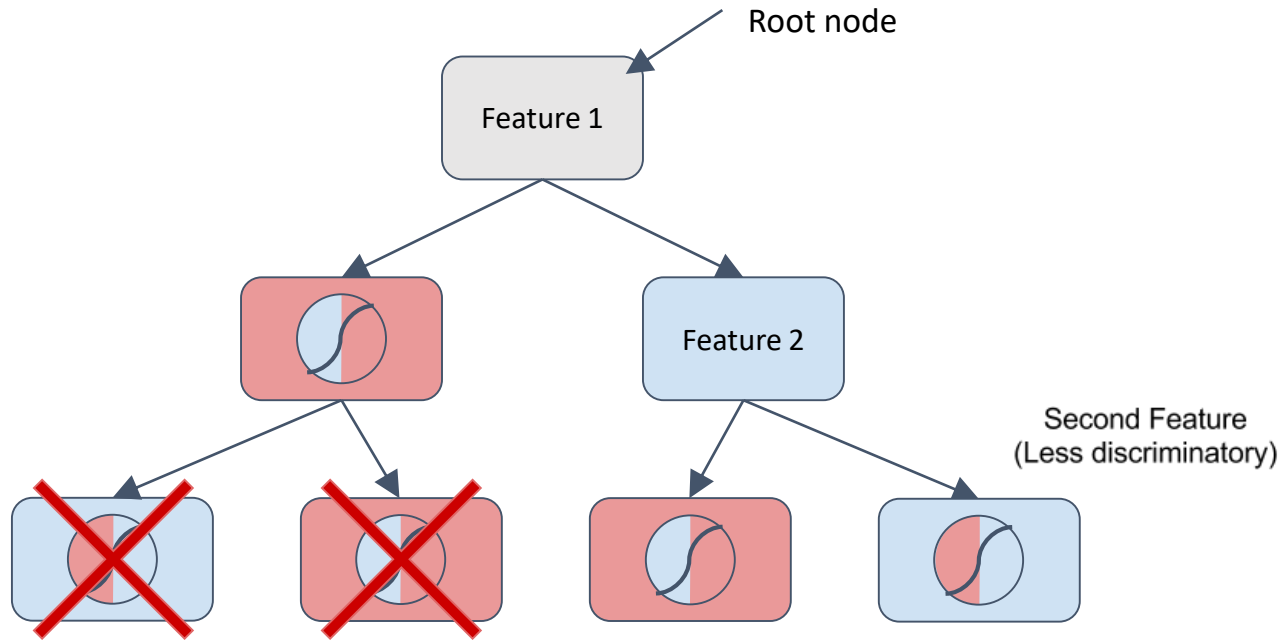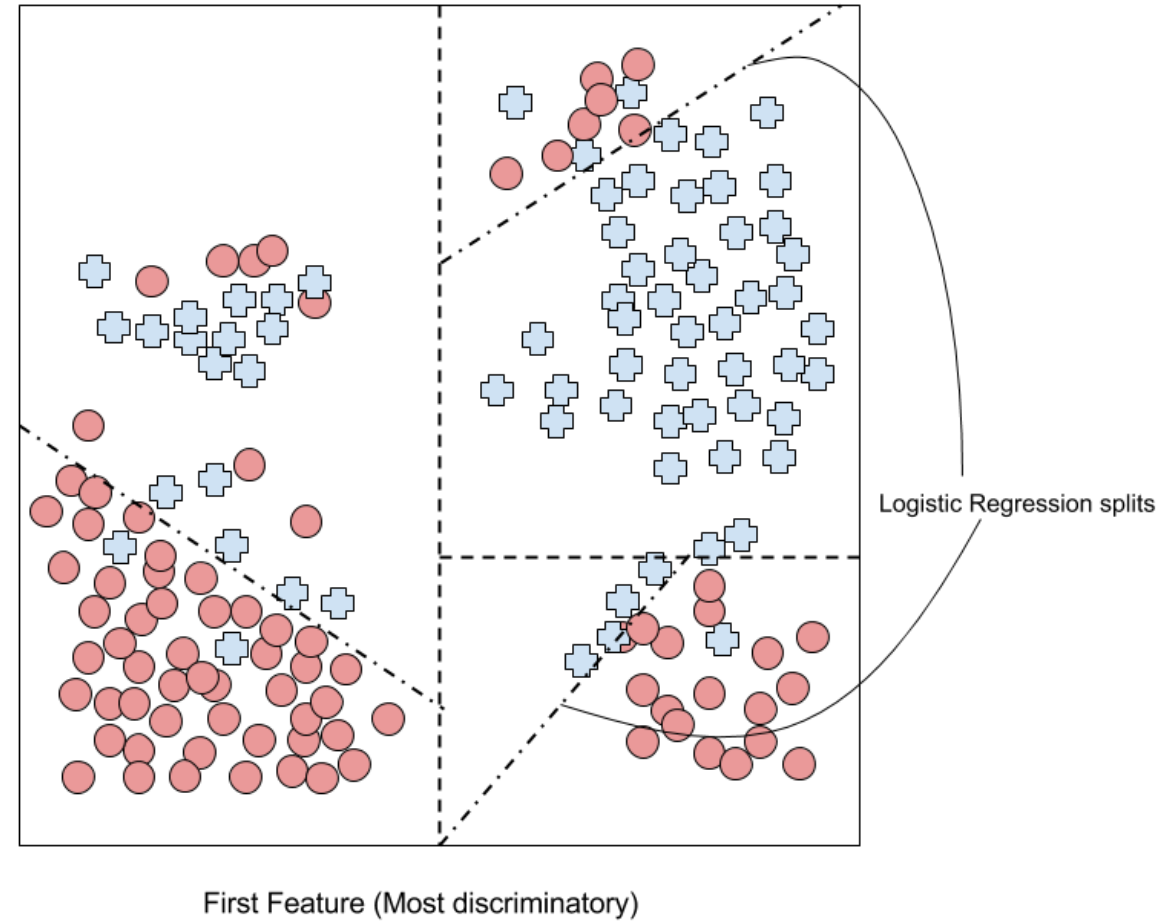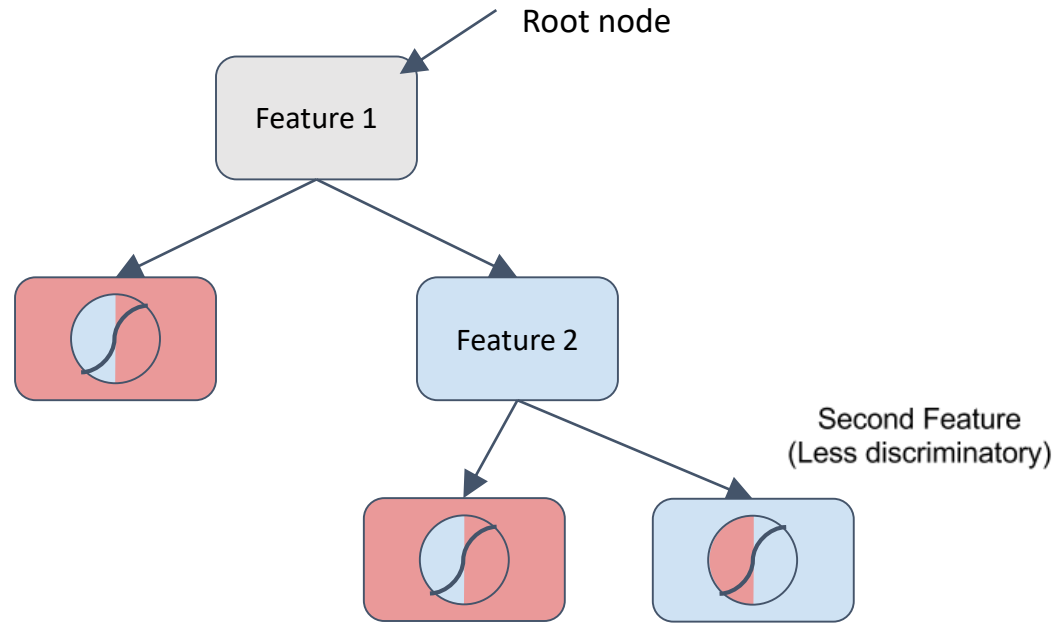(Less discriminatory)

First Feature (Most discriminatory)
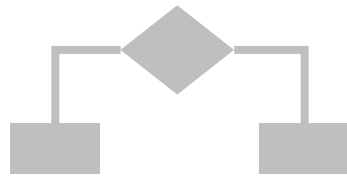
Logistic Regression splits

# LOGISTIC MODEL TREES

**DISTRIBUTED IMPLEMENTATION**

**Spark's Decision Tree**
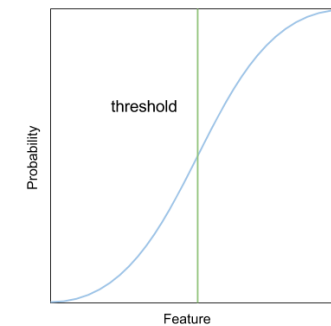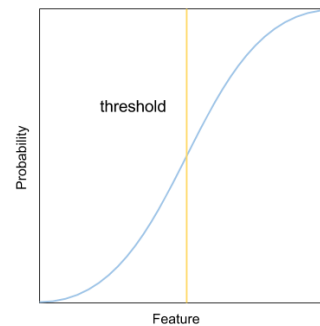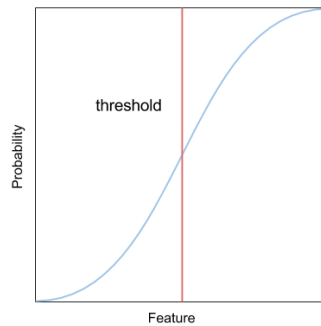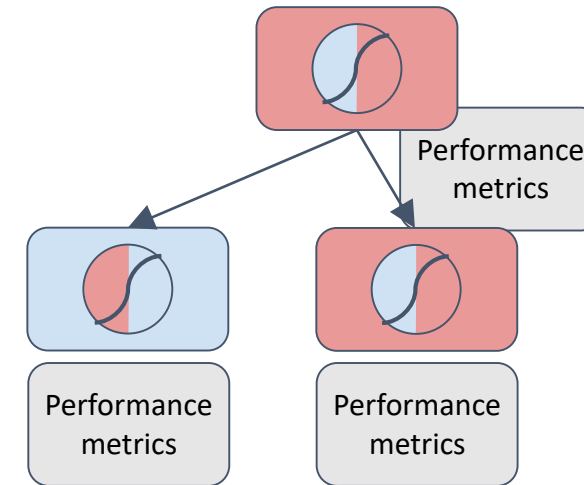(distributed implementation of random forests)

**Spark's Logistic Regression / weka's
Logistic Regression on the nodes**

LMT Cost function to fix the logistic regression threshold

- AccuracyCostFunction
- ConfusionMatrix
- PrecisionCostFunction
- PrecisionRecallCostFunction
- RocCostFunction

The same cost function for pruning criteria

## ADVANTAGES OF THIS IMPLEMENTATION

### Big datasets
Power of spark to distribute building the tree and logistic regressions

### Medium datasets
Distributed tree growth and weka's logistic regression

### Small datasets
Although it can be slow to distribute the data for the decision tree, cost functions can be still used and specific optimization for particular cases

# Example of DLMT algorithm in a synthetic dataset

# AUTOMATED BENCHMARKING FRAMEWORK

- Metrics
- Demo

@StratioBD

|  |  | PREDICTION | |
|---|---|---|---|
|  |  | Positive | Negative |
| **TRUE CONDITION** | Positive | True Positives | False Negatives |
|  | Negative | False Positives | True Negatives |

→ **True Positive Rate** (Recall)

→ **False Positive Rate**

↓ **Precission**

**TPR** = TP/(TP+FN) Insensitive to unbalance

**FPR** = FP/(FP+TN) Insensitive to unbalance
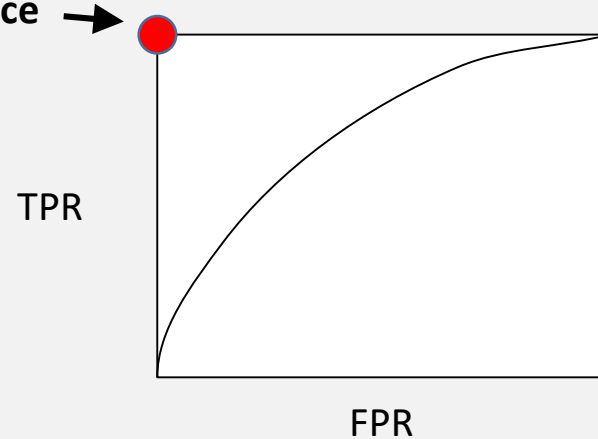
**Precision** = TP/(TP+FP) Sensitive to unbalance

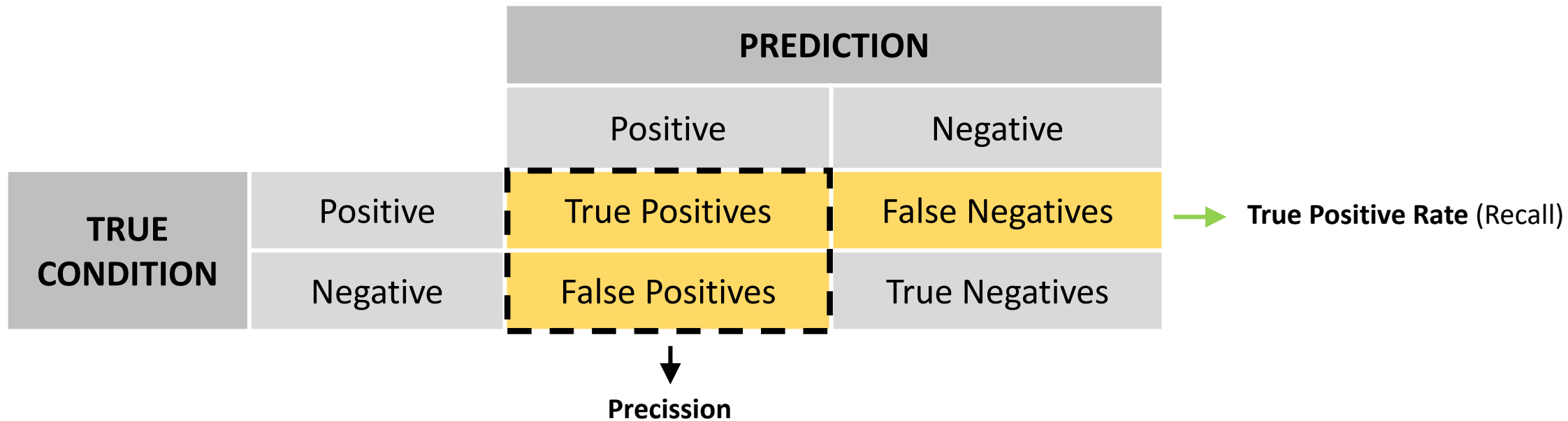**Accuracy** = (TP+TN)/(TP+TN+FP+FN) Sensitive to unbalance

|  | | PREDICTION | |
|---|---|---|---|
|  | | Positive | Negative |
| **TRUE CONDITION** | Positive | True Positives | False Negatives |
|  | Negative | False Positives | True Negatives |

→ **True Positive Rate** (Recall)

→ **False Positive Rate**

**AUROC (AUC): TPR/FPR** -> Insensitive to unbalance!

**Best performance** →

TPR

FPR

**PREDICTION**

| | | Positive | Negative |
|---|---|---|---|
| **TRUE CONDITION** | Positive | True Positives | False Negatives | → **True Positive Rate** (Recall) |
| | Negative | False Positives | True Negatives |

↓

**Precission**

**AUPRC: Precision/TPR** -> Sensitive to unbalance!

← **Best performance**

Precision

Recall

Data

Algorithms $\begin{bmatrix} f_1 \\ \cdots \\ f_n \end{bmatrix}$

*ABF*

**Benchmark**

# BENCHMARKING RESULTS

**1** Accuracy **VS** Explainability

**2**

**3** **Performance Metrics:**
AUROC, AUPRC, ACCURACY

**4** **Automatic Benchmarking Framework**

$f_1$
$f_n$

*ABF*

Benchmark

# THANK YOU

**UNITED STATES**

Tel: (+1) 408 5998830

**EUROPE**

Tel: (+34) 91 828 64 73

contact@stratio.com

www.stratio.com

@StratioBD

**STRATIO**®

# WE ARE HIRING

people@stratio.com

@StratioBD