

Selecting Your Algorithm



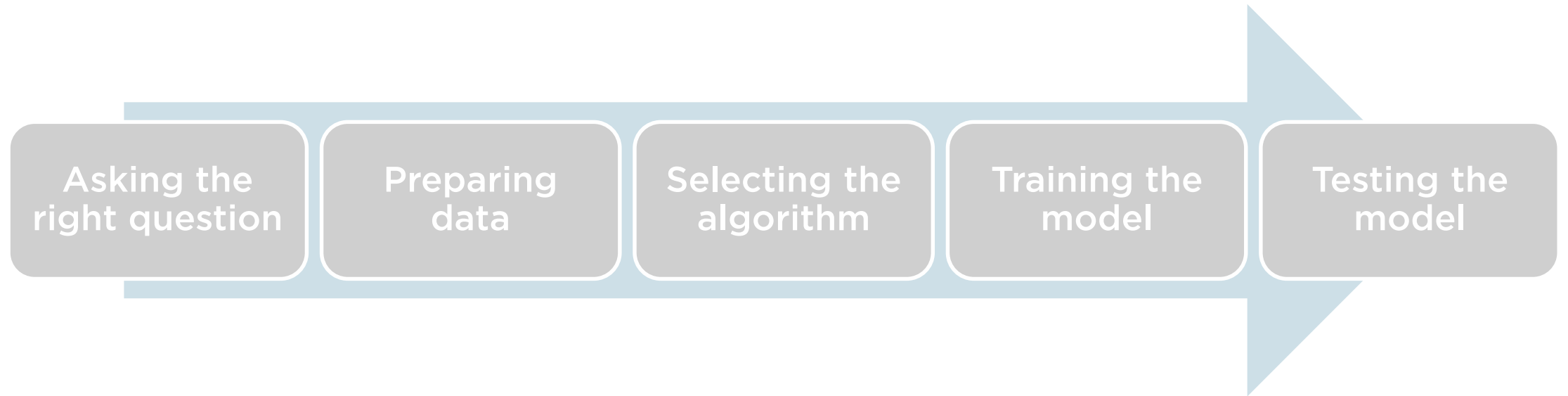
Jerry Kurata

CONSULTANT

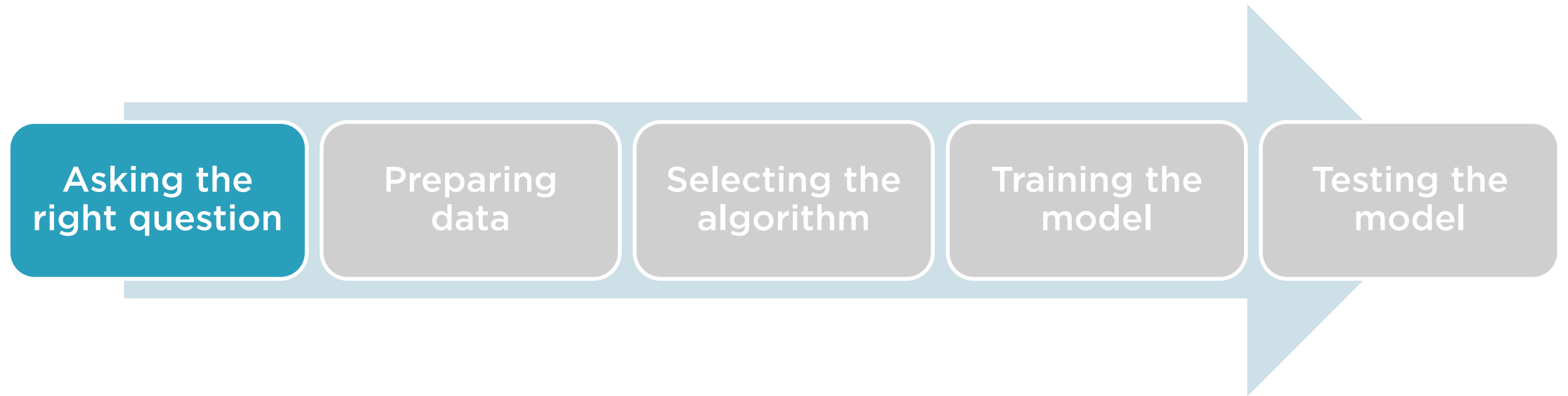
@jerrykur www.insteptech.com



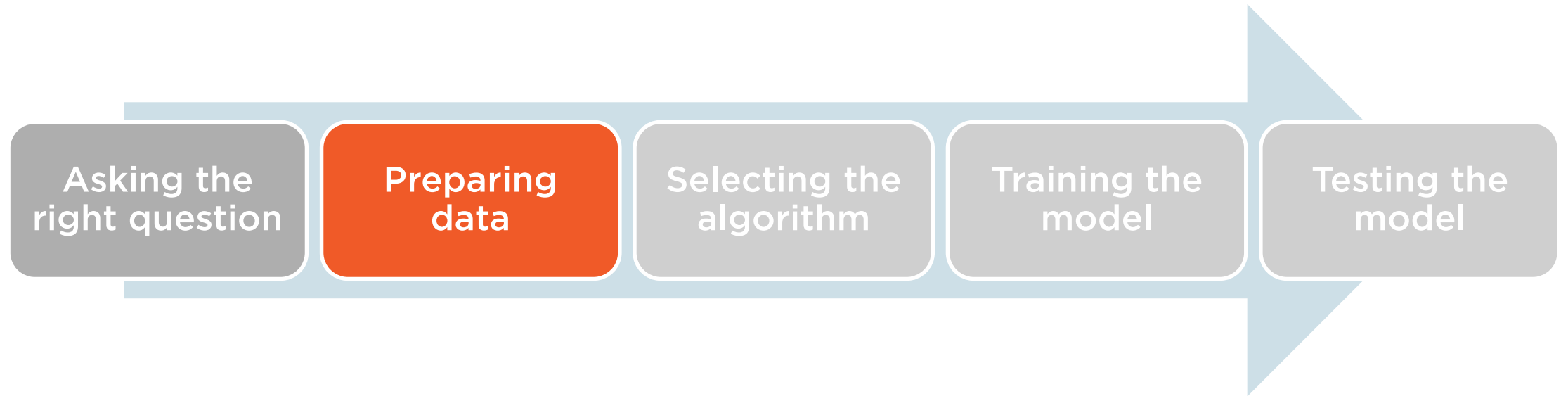
Machine Learning Workflow



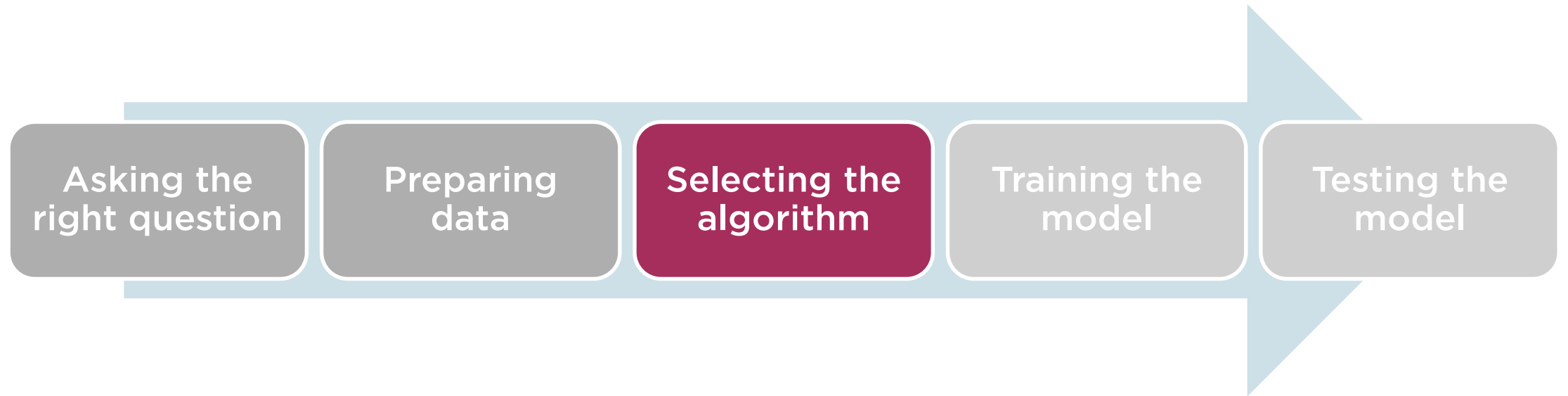
Machine Learning Workflow



Machine Learning Workflow



Machine Learning Workflow



Overview



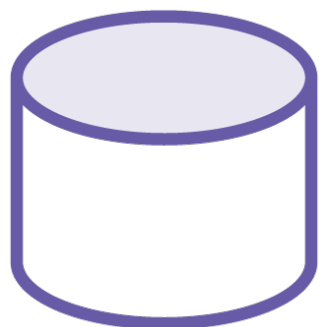
Role of algorithm

Perform algorithm selection

- Use solution statement to filter algorithms
- Discuss best algorithms
- Select one initial algorithm

Role of Algorithm

train()



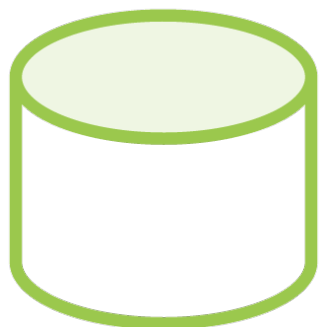
Training
Data



Algorithm



predict()



Real
Data



Model



A = 2.2

B = 0.4

C = 9.2



Over 50 algorithms



Algorithm Selection

Compare factors

**Difference of opinions about which factors
are important**

You will develop your own factors



Algorithm Decision Factors

Learning Type

Result

Complexity

Basic vs Enhanced



Learning Type



Learning Type

“Use the Machine Learning Workflow to process and transform DOT data to create a prediction model. This model must predict whether a flight would arrive 15+ minutes after the scheduled arrival time with 70+% accuracy.”



Learning Type

*“Use the Machine Learning Workflow to process and transform DOT data to create a **prediction model**. This model must predict whether a flight would arrive 15+ minutes after the scheduled arrival time with 70+% accuracy.”*

Prediction Model => Supervised machine learning



Over ~~50~~ 28 algorithms



Result Type

Regression

- Continuous values
- $\text{price} = A * \# \text{ bedroom} + B * \text{size} + \dots$

Classification

- Discrete values
- small, medium, large
- 1-100, 101-200, 201-300
- true or false



Result Type

“... predict whether a flight would arrive 15+ minutes after the scheduled arrival time .”



Result Type

“... predict whether a flight would arrive 15+ minutes after the scheduled arrival time .”

ARR_DEL15

Binary (TRUE/FALSE)

Algorithm must support classification

- Binary classification



Over ~~50~~ ~~28~~ 20 algorithms



Complexity

Keep it Simple

Eliminate “ensemble” algorithms

- Container algorithm
- Multiple child algorithms
- Boost performance
- Can be difficult to debug



Over ~~50~~ ~~28~~ ~~20~~ 14 algorithms



Enhanced vs. Basic

Enhanced

- Variation of Basic
- Performance improvements
- Additional functionality
- More complex

Basic

- Simpler
- Easier to understand



Candidate Algorithms

Naive Bayes

**Logistic
Regression**

Decision Tree

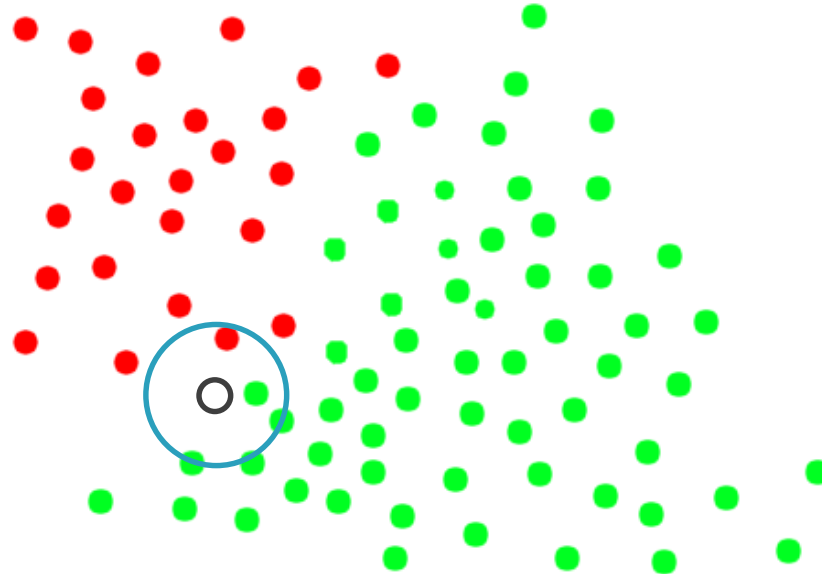


Naive Bayes

Based on likelihood
and probability

Every feature has the
same weight

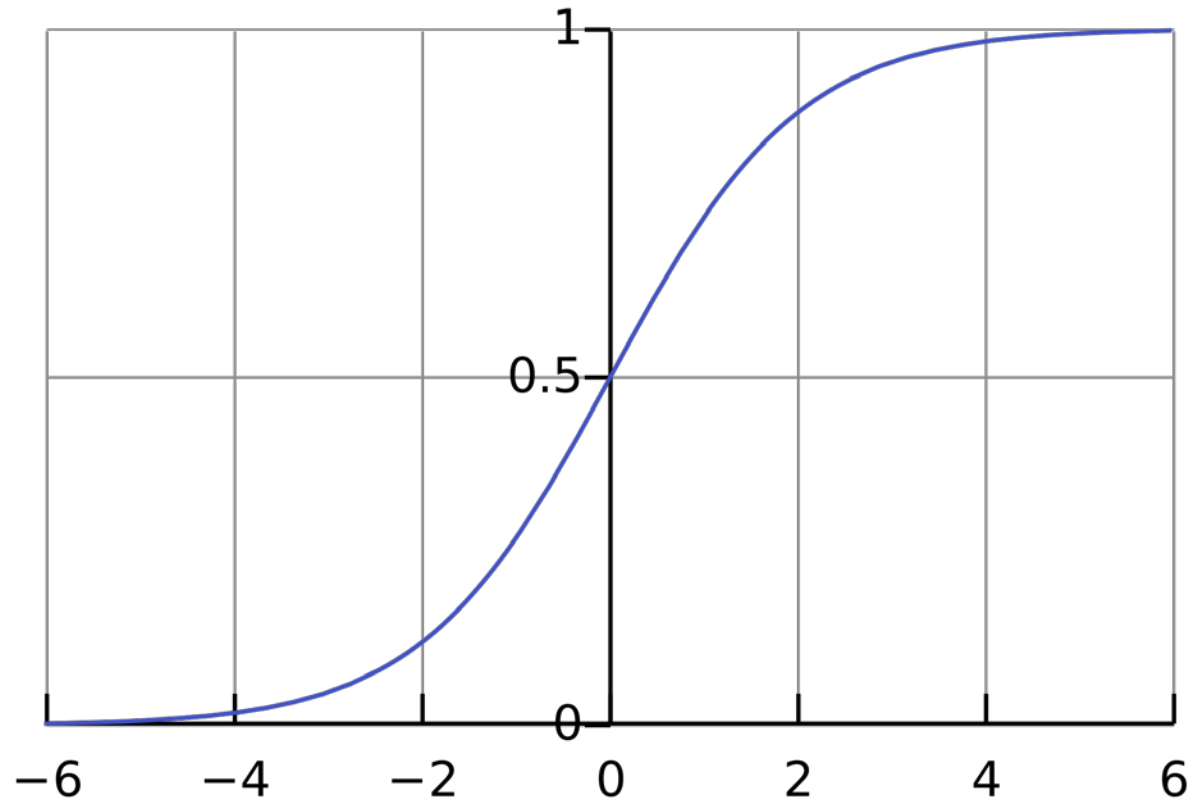
Requires smaller
amount of data



Logistic Regression

Confusing name,
binary result

Relationship between
features are weighted



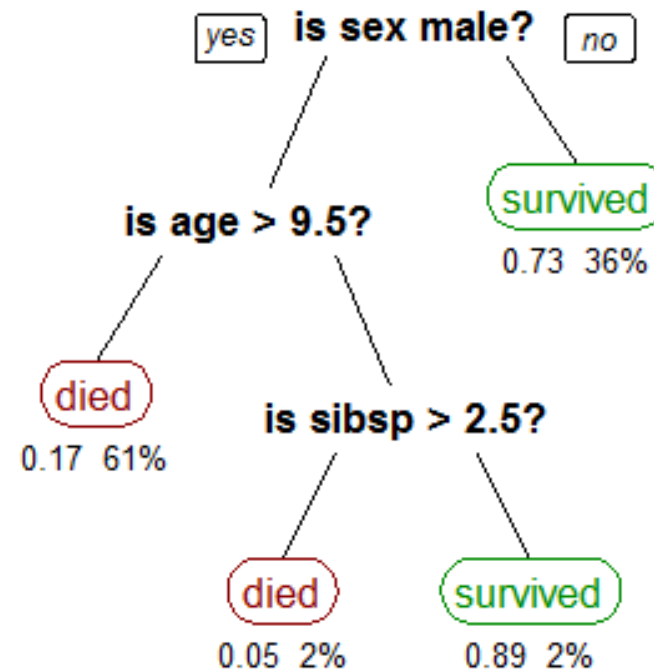
Decision Tree

Binary Tree

Node contains decision

Requires enough data to determine nodes and splits

Titanic Survival



Selected Algorithm

Logistic
Regression

Simple - easy to understand

Fast - up to 100X faster

Stable to data changes



Summary



Lots of algorithms available

Selection based on

- Learning = Supervised
- Result = Binary classification
- Non-ensemble
- Basic

Logistic Regression selected for training

- Simple, fast, and stable

