

# Data Science Basics and Decision Tree Models

# What Really *is* Data Mining (or Data Science) ?

- A process for using information technology to extract useful (non-trivial, hopefully actionable) knowledge from large bodies of data

# Data Science Tasks

Many problems have one of these data science tasks:

- Classification
  - Predict for each individual in a population which class this individual belongs to.
- Regression
  - “value estimation”, estimate the numerical value of some variable specific to an individual.
- Similarity matching
  - Identify similar individuals based on data known about them.
- Clustering
  - Group individuals by their similarity.

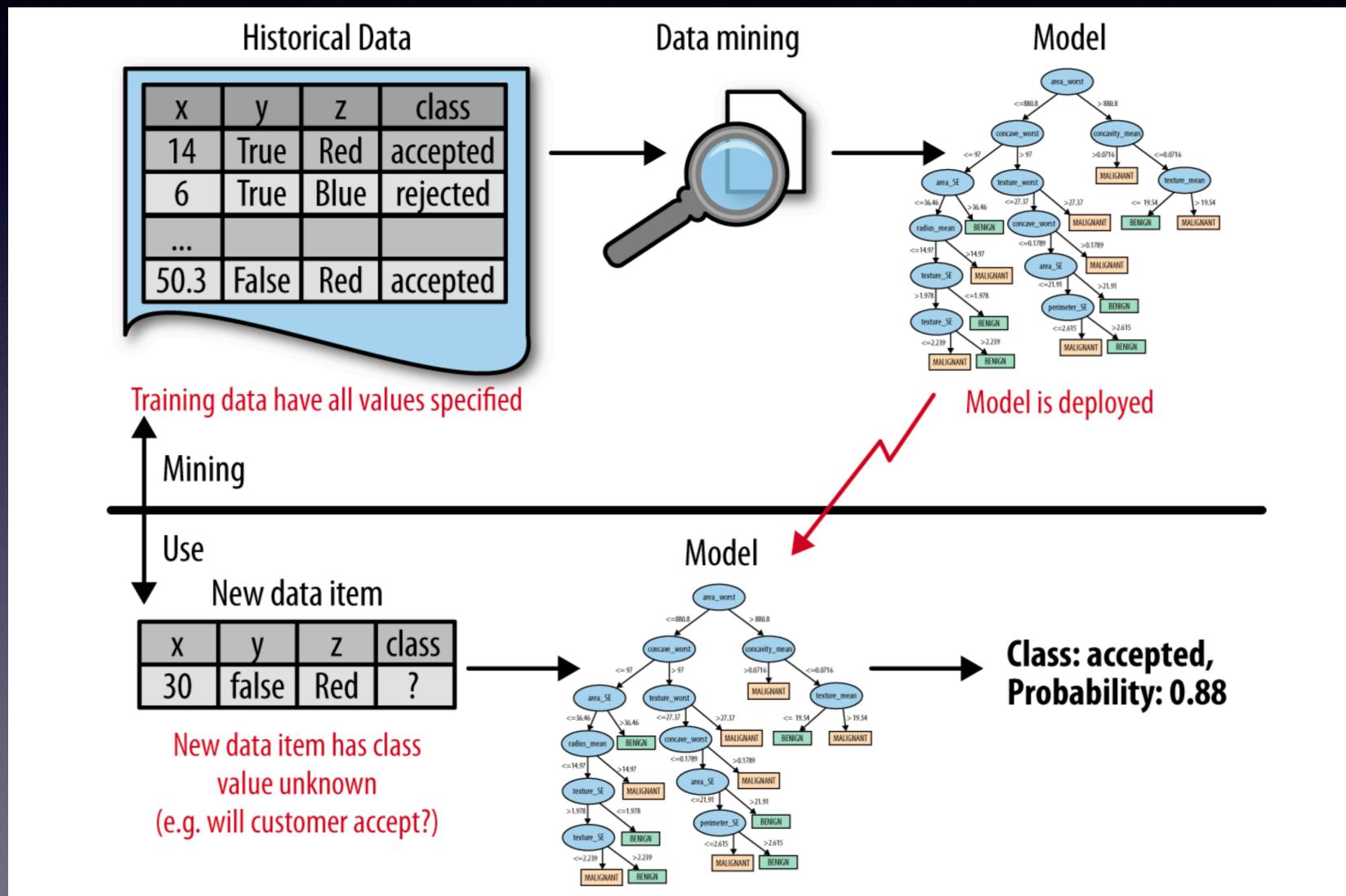
# Data Science Tasks

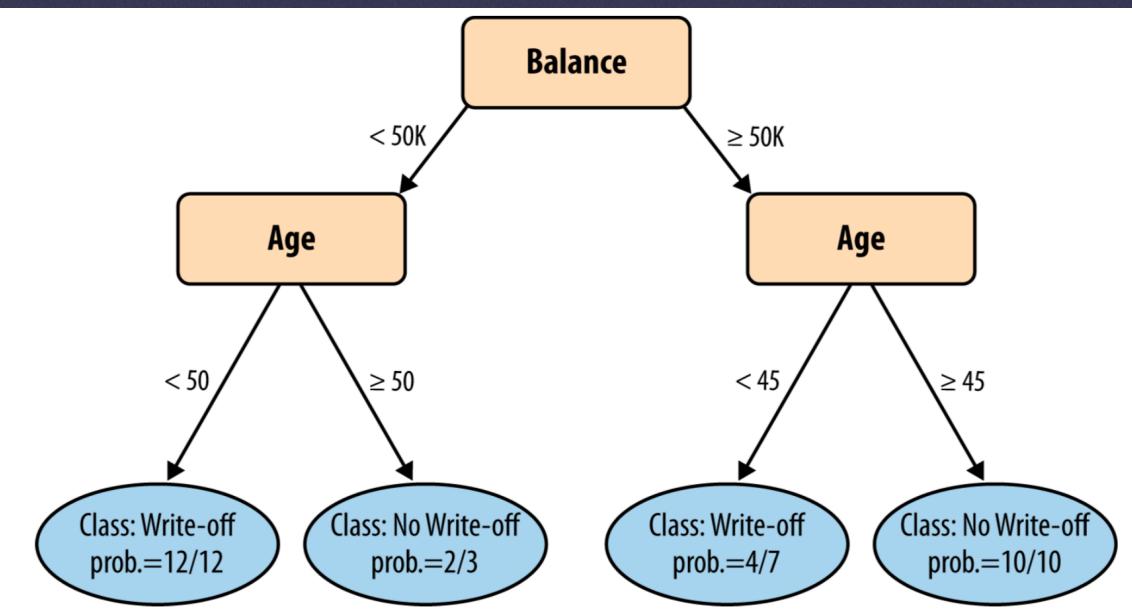
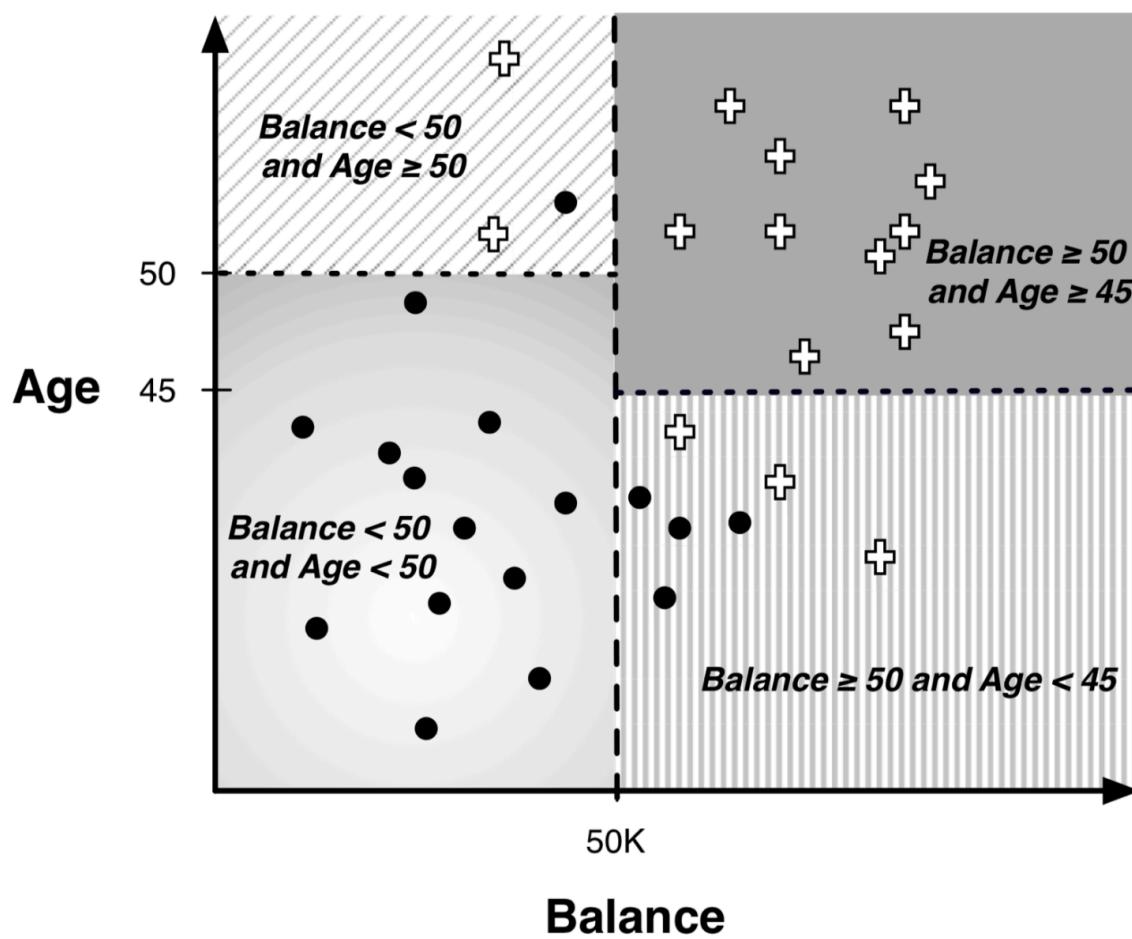
- Co-occurrence grouping
  - Find associations between entities based on transactions involving them. The result is a description of items that occur together.
- Profiling
  - “behavior description”, characterize the typical behavior of an individual, group, or population.
- Link prediction
  - Predict connections between data items, usually. By suggesting that a link should exist, and possibly also estimating the strength of the link.

# Supervised vs. Unsupervised Learning

- Supervised methods
  - Classification, regression
- Unsupervised methods
  - Clustering, profiling, co-occurrence grouping
- Technically, there must be data on the target for supervised learning methods.
  - Classification and regression are distinguished by the type of target. Regression involve a numeric target while classification involves a categorical target.

# Classification Example





# Example dataset for a Supervised Classification Problem

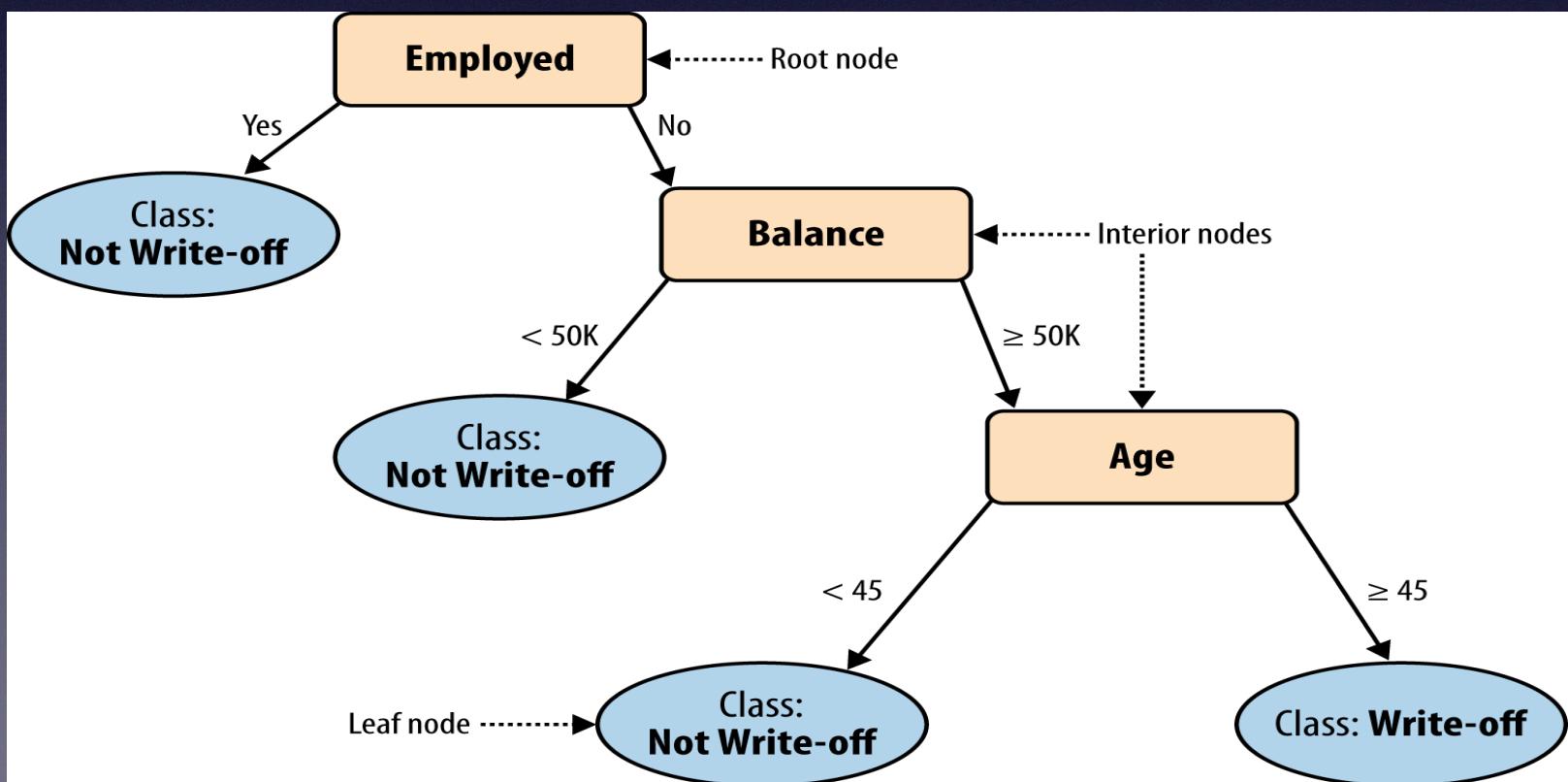
The diagram illustrates a dataset for a supervised classification problem. At the top, a bracket labeled "Attributes" spans across the first four columns of the table. A separate bracket labeled "Target attribute" points to the last column, "Write-off". An arrow on the left side of the table points to the fifth row, highlighting it as an example.

Name	Balance	Age	Employed	Write-off
Mike	\$200,000	42	no	yes
Mary	\$35,000	33	yes	no
Claudio	\$115,000	40	no	no
Robert	\$29,000	23	yes	yes
Dora	\$72,000	31	no	no

This is one row (example).  
Feature vector is: <Claudio,115000,40,no>  
Class label (value of Target attribute) is no

# Tree-Structured Models

- Classify Claudio
  - Balance=115K, Employed=No, and Age=40



# Tree-Structured Models: “Rules”

- No two parents share descendants
- There are no cycles
- The branches always “point downwards”
- Every example always ends up at a leaf node with some specific class determination

# Tree Induction

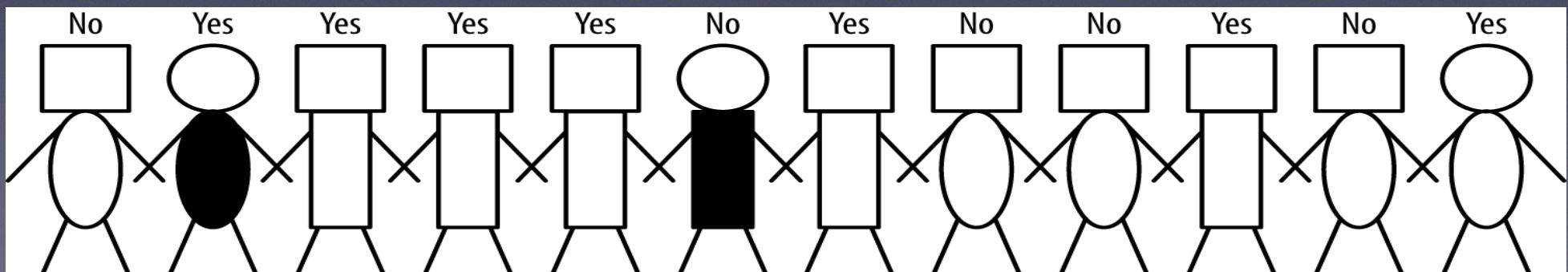
- How do we create a classification tree from data?
  - **Divide-and-conquer** approach
  - Take each data subset and *recursively* apply attribute selection to find the best attribute to partition it (create the “purest” subgroups possible using the attributes)
- When do we stop?
  - The nodes are pure
  - There are no more variables (don’t use the same variable more than once), or even earlier (over-fitting)

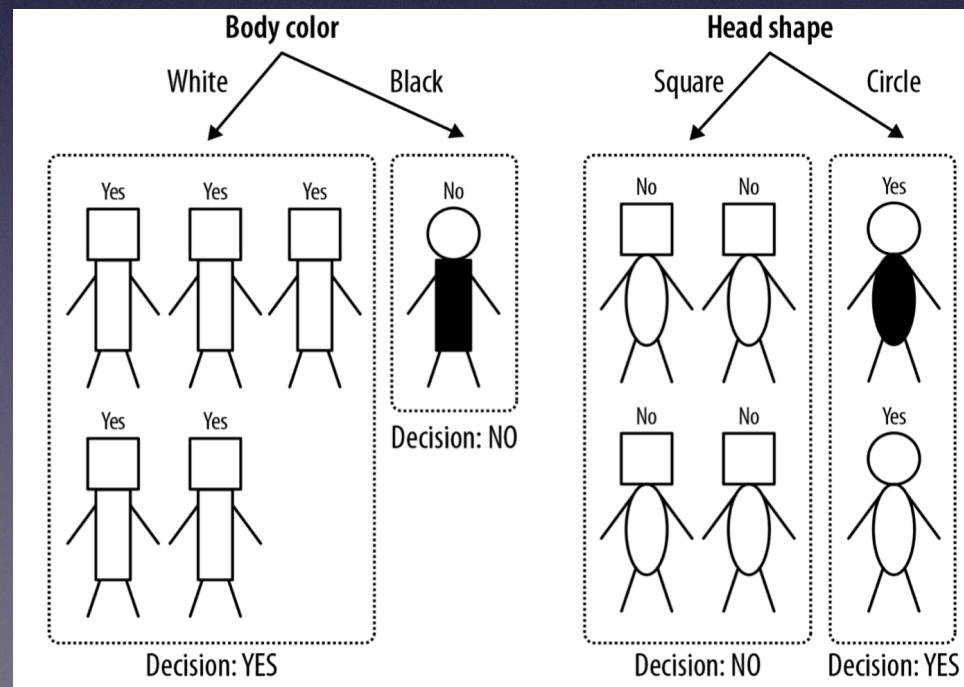
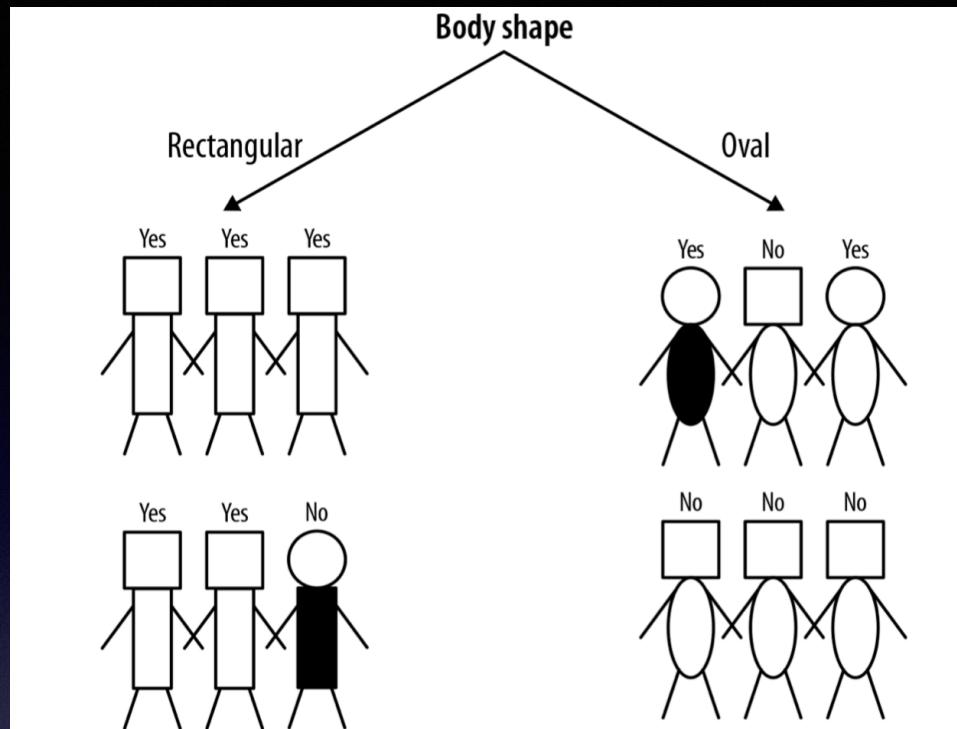
# What are the Questions?

- How can we segment the population into groups that differ from each with respect to some quantity of interest?
- How can we judge whether a variable contains important information about the target variable?
  - How much?

# Create a classification tree from data

- Objective: Based on customer attributes, partition the customers into subgroups that are less impure – with respect to the class (i.e., such that in each group as many instances as possible belong to the same class)





# Calculating Impurity (as Entropy)

- Most common splitting criterion: information gain
- Entropy: impurity measure

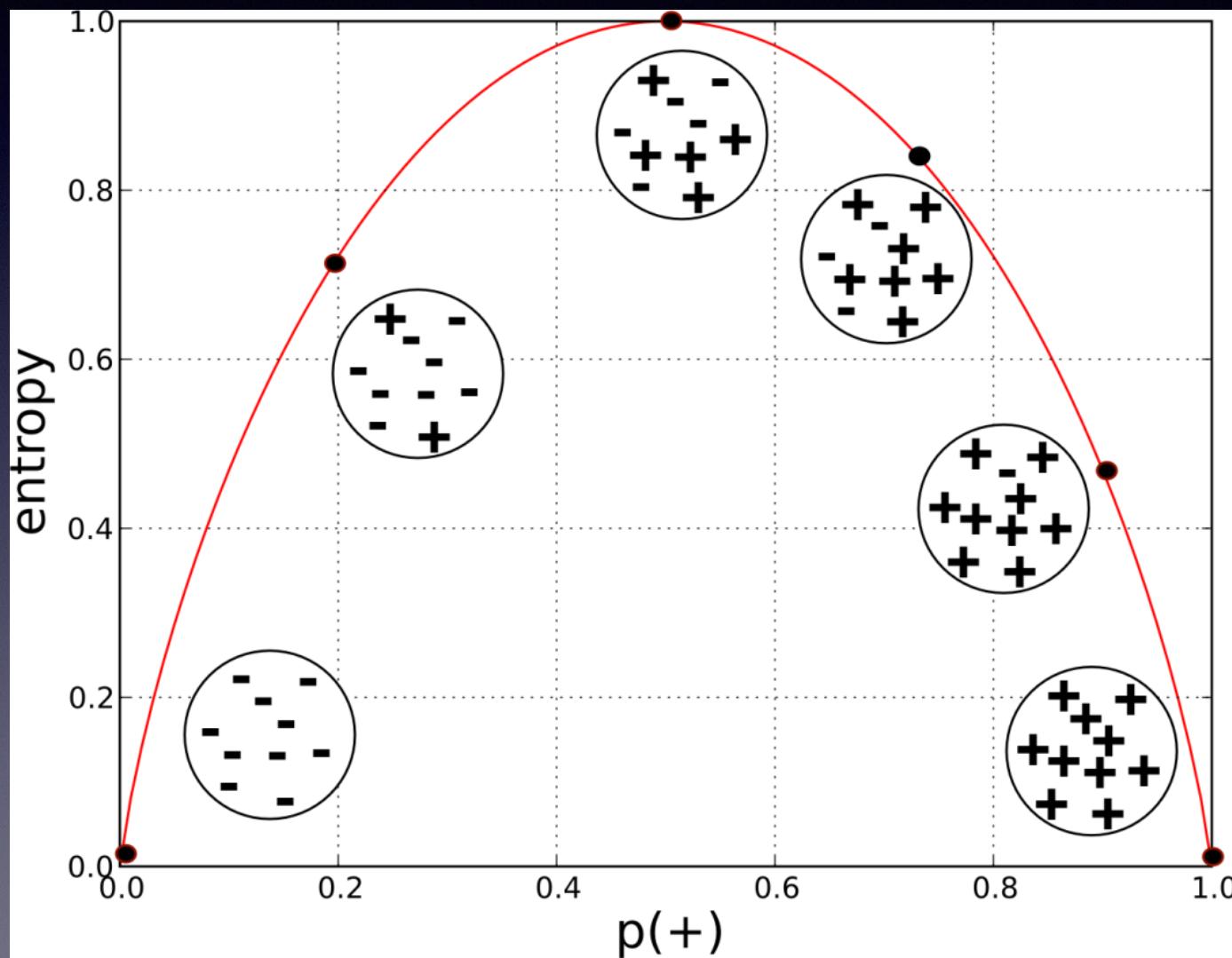
$$\text{entropy} = - p_1 \log(p_1) - p_2 \log(p_2) - \dots$$

$p_i$  is the proportion of class  $i$  in the data

- Example: the population is composed of 14 cases of class “Yes” and 16 cases of class “No”

$$-\left(\frac{14}{30} \cdot \log_2 \frac{14}{30}\right) - \left(\frac{16}{30} \cdot \log_2 \frac{16}{30}\right) = 0.996$$

# Entropy Function Graph of a Two-class Dataset

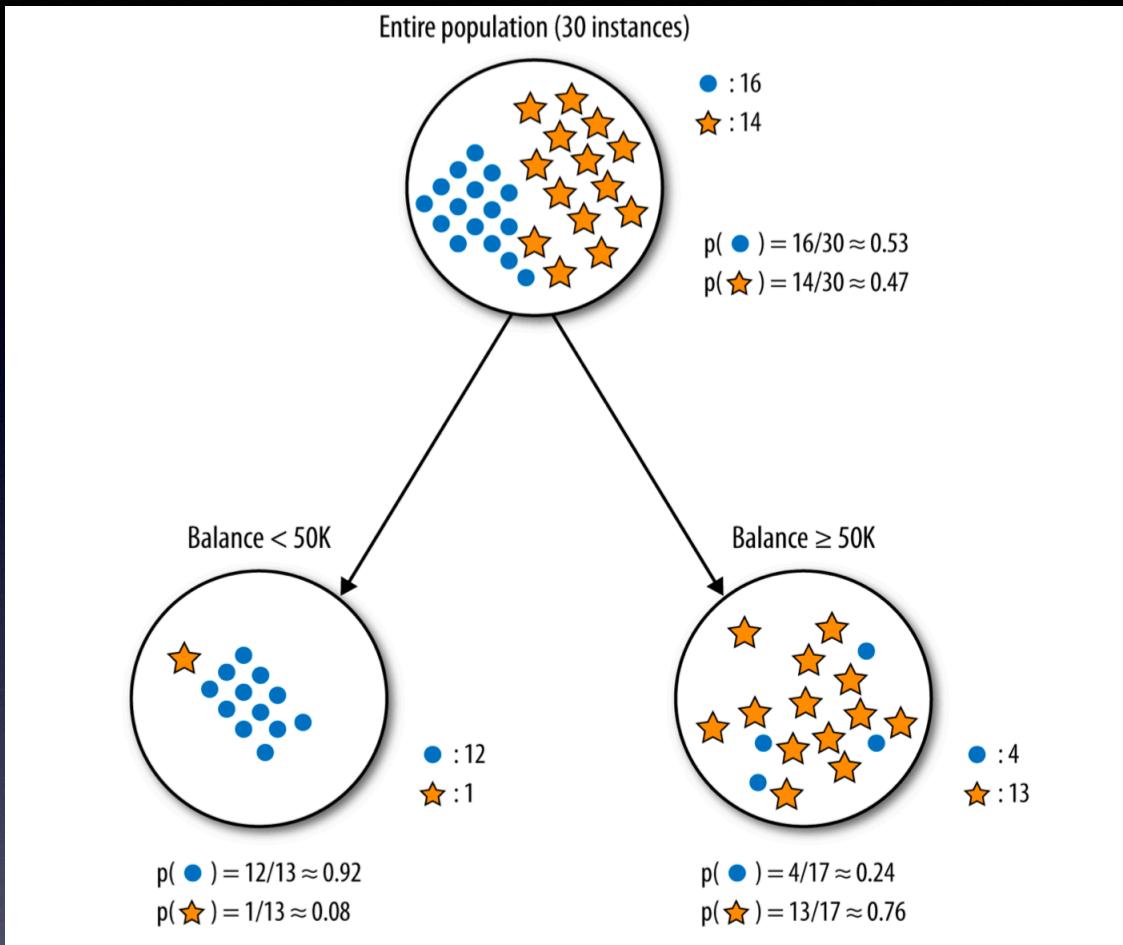


# Information Gain

- Information Gain (IG): measure how much an attribute improves (decreases) entropy over the whole segmentation it creates (measures the change in entropy)

$$IG(\text{parent}, \text{children}) = \text{entropy}(\text{parent}) - [p(c_1) \times \text{entropy}(c_1) + p(c_2) \times \text{entropy}(c_2) + \dots]$$

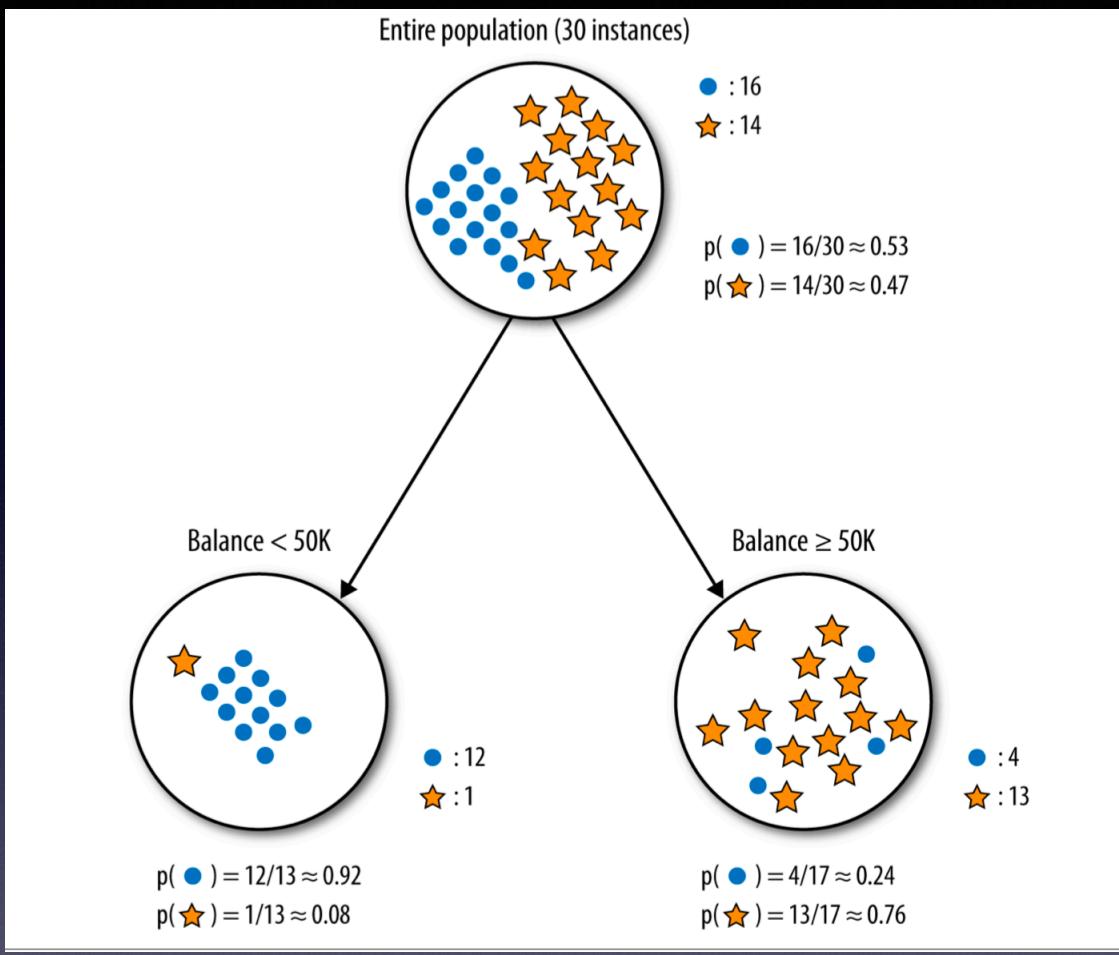
- The entropy for each child ( $c_i$ ) is weighted by the proportion of instances belonging to that child,  $p(c_i)$ .



$$entropy = - p_1 \log(p_1) - p_2 \log(p_2) - \dots$$

- The information gain of this split is?

$$\begin{aligned}
 \text{entropy}(\text{parent}) &= -[p(\bullet) \times \log_2 p(\bullet) + p(\star) \times \log_2 p(\star)] \\
 &\approx -[0.53 \times -0.9 + 0.47 \times -1.1] \\
 &\approx 0.99 \quad (\text{very impure})
 \end{aligned}$$



$$\begin{aligned}
 \text{entropy}(\text{Balance} < 50K) &= -[p(\bullet) \times \log_2 p(\bullet) + p(\star) \times \log_2 p(\star)] \\
 &\approx -[0.92 \times (-0.12) + 0.08 \times (-3.7)] \\
 &\approx 0.39
 \end{aligned}$$

$$\begin{aligned}
 \text{entropy}(\text{Balance} \geq 50K) &= -[p(\bullet) \times \log_2 p(\bullet) + p(\star) \times \log_2 p(\star)] \\
 &\approx -[0.24 \times (-2.1) + 0.76 \times (-0.39)] \\
 &\approx 0.79
 \end{aligned}$$

- The information gain of this split is:

$$\begin{aligned}
 \text{IG} &= \text{entropy}(\text{parent}) - [p(\text{Balance} < 50K) \times \text{entropy}(\text{Balance} < 50K) \\
 &\quad + p(\text{Balance} \geq 50K) \times \text{entropy}(\text{Balance} \geq 50K)] \\
 &\approx 0.99 - [0.43 \times 0.39 + 0.57 \times 0.79] \\
 &\approx 0.37
 \end{aligned}$$

# Why trees?

- Decision trees (DTs), or classification trees, are one of the most popular data mining tools
- They are:
  - Easy to understand and explain
  - Easy to implement
  - Easy to use
  - Computationally cheap
- Almost all data mining packages include DTs

# Attribute Selection

- Reasons for selecting only a subset of attributes:
  - Better insights and business understanding
  - Better explanations and more tractable models
  - Reduced cost
  - Faster predictions
  - Better predictions!
    - Too many attributes bring the over-fitting problem
  - And also determining the most informative attributes.

# Hands-on: Decision Tree

- Download Class 4 lab from blackboard