CS57300
PURDUE UNIVERSITY
JANUARY 22, 2019

# DATA MINING

# HYPOTHESIS TESTING

# TYPES OF HYPOTHESES

Broad categories

▸ **Descriptive**: propositions that describe a characteristic of an object

▸ **Relational**: propositions that describe relationship between 2+ variables

▸ **Causal**: propositions that describe the effect of one variable on another

Specific characteristics

▸ **Non-directional**: an differential outcome is anticipated but the specific nature of it is not known (e.g., the tuning parameter will affect algorithm performance)

▸ **Directional**: a specific outcome is anticipated (e.g., the use of pruning will increase accuracy of models compared to no pruning)

**Descriptive Hypothesis**

**Non–Directional Relational Hypothesis**

**Directional Relational Hypothesis**

**Directional Causal Hypothesis**

Stronger

# HYPOTHESES EXAMPLE

▸ The query response time is measured for a few different search engines

▸ Different hypotheses

  ▸ **Descriptive:** The query response time for Google follows a normal distribution

  ▸ **Non-directional relational:** The average response time for a new search engine, QuickSearch, is different from Google's average response time

  ▸ **Directional relational:** The average response time of QuickSearch is shorter than that of Google's

  ▸ **Directional causal:** The response time of QuickSearch is shorter than Google's because they cache results of more queries
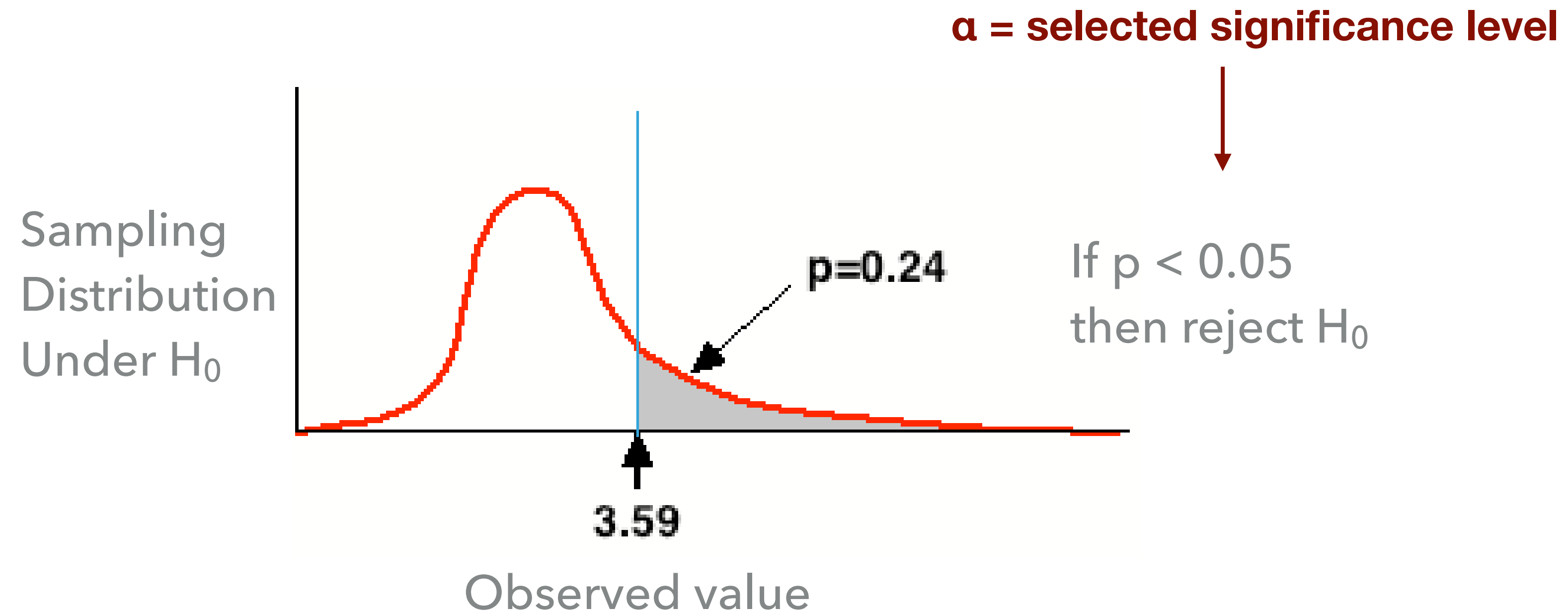
# HYPOTHESIS TESTING

▸ Statistical hypothesis test is a method used in statistics that tells you the likelihood of a specific result would happen by chance

▸ **Null hypothesis** ($H_0$):

  ▸ Presumed true until statistical inference indicates otherwise; set up to be refuted by alternative

▸ **Alternative hypothesis** ($H_1$):

  ▸ Rival hypothesis; that we conjecture is true

▸ Assuming the null hypothesis is true, what's the probability of getting a statistic that is at least as extreme as the statistic that was actually obtained through the data?

# HYPOTHESIS TESTING STRATEGY

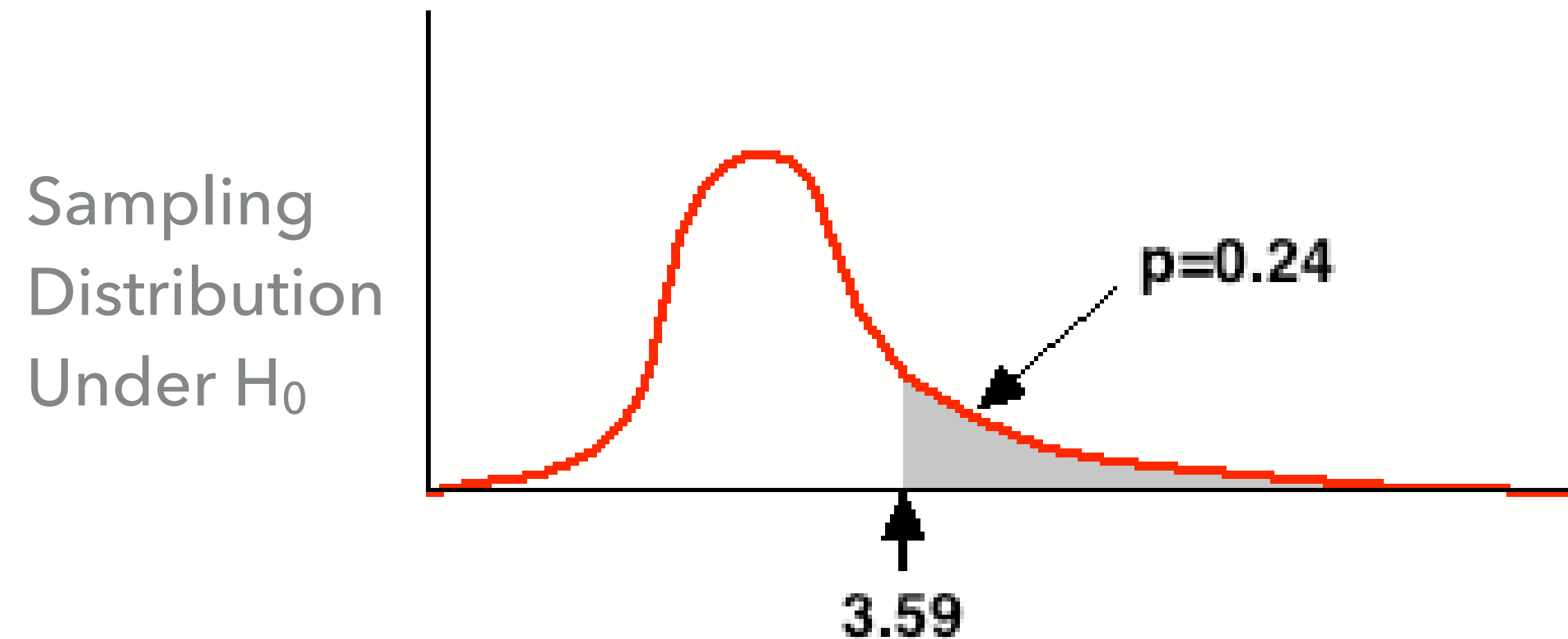https://towardsdatascience.com/data-science-simplified-hypothesis-testing-56e180ef2f71

▸ Formulate null and alternative hypothesis

   ▸ $H_0$: QuickSearch' mean response time = Google's mean response time

   ▸ $H_1$: QuickSearch' mean response time ≠ Google's mean response time

▸ Gather a sample statistic (e.g., $\delta$ =difference of QuickSearch's and Google's mean response time)

▸ Determine the sampling distribution for the statistic under the null hypothesis

▸ Use the sampling distribution to calculate the probability of obtaining the observed value of $\delta$, given $H_0$

   ▸ If the probability is low, reject $H_0$ in favor of $H_1$

# REJECTING THE NULL HYPOTHESIS

**α = selected significance level**

Sampling
Distribution
Under $H_0$

p=0.24

If p < 0.05
then reject $H_0$

3.59

Observed value

# STATISTICAL SIGNIFICANCE

▸ A value of a statistic is **statistically significant** if it is unlikely to occur under the null hypothesis

Sampling
Distribution
Under $H_0$

p=0.24

3.59

**significance level**  $\alpha = p(reject\ H_0 | H_0\ true) = p(type\ 1\ error)$

# ERRORS

| Truth | | Decision | |
|---|---|---|---|
| | | Reject $H_0$ | Don't reject $H_0$ |
| Truth | $H_0$ | *Type 1 error* | |
| | $H_1$ | | *Type 2 error* |

▸ Type 1: null is rejected when it is true

   ▸ E.g., conclude cancer drug increases life expectancy when in fact it doesn't

   ▸ Generally considered to be most serious error

▸ Type 2: null is accepted when it is false

   ▸ E.g., conclude that cancer drug does not increase life expectancy when in fact it does

# STATISTICAL POWER

▸ Lack of statistical significance does not necessarily imply that $H_0$ is true

▸ Test could have low statistical power: $(1 - \beta)$ **portion of sampling distribution for alternative that is above threshold**



$$\beta = p(accept\ H_0 | H_0\ false) = p(type\ 2\ error)$$

# HOW TO INCREASE POWER
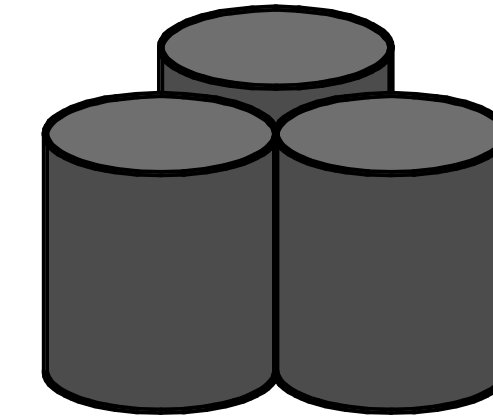
▸ Increase sample size

▸ Decrease sample variability

   ▸ Matching, sample selection, control for confounding variables, increase precision of measurements

▸ Increase effect size

   ▸ More extreme experimental conditions, avoid ceiling/floor effects

▸ Increase alpha (e.g., from 0.05 to 0.10, but this increases type 1 errors)
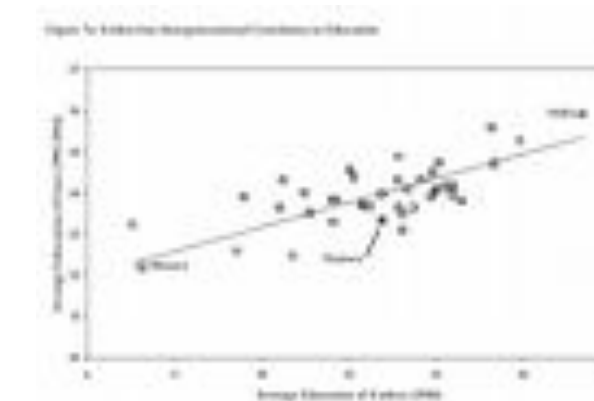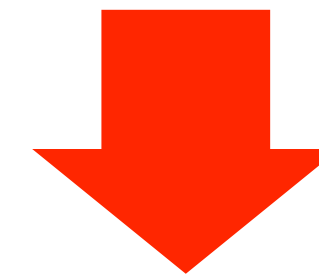
# DATA AND MEASUREMENT

# REFLECTING REAL WORLD THROUGH DATA



Real world

Data

Relationship
in real world

Relationship
in data

Goal: map domain entities to symbolic representations

# WHAT IS DATA?

▸ Collection of entities and their attributes

▸ **Attribute**: property or characteristic of an entity (e.g., eye color, temperature)

▸ **Entity**: collection of attributes
Aka: record, point, case, sample, object, or instance

**Attributes**

**Entities**

| Name | Thread pitch (mm) | Minor diameter tolerance | Nominal diameter (mm) | Head shape | Price for 50 screws | Available at factory outlet? | Number in stock | Flat or Phillips head? |
|---|---|---|---|---|---|---|---|---|
| M4 | 0.7 | 4g | 4 | Pan | $10.08 | Yes | 276 | Flat |
| M5 | 0.8 | 4g | 5 | Round | $13.89 | Yes | 183 | Both |
| M6 | 1 | 5g | 6 | Button | $10.42 | Yes | 1043 | Flat |
| M8 | 1.25 | 5g | 8 | Pan | $11.98 | No | 298 | Phillips |
| M10 | 1.5 | 6g | 10 | Round | $16.74 | Yes | 488 | Phillips |
| M12 | 1.75 | 7g | 12 | Pan | $18.26 | No | 998 | Flat |
| M14 | 2 | 7g | 14 | Round | $21.19 | No | 235 | Phillips |
| M16 | 2 | 8g | 16 | Button | $23.57 | Yes | 292 | Both |
| M18 | 2.1 | 8g | 18 | Button | $25.87 | No | 664 | Both |
| M20 | 2.4 | 8g | 20 | Pan | $29.09 | Yes | 486 | Both |
| M24 | 2.55 | 9g | 24 | Round | $33.01 | Yes | 982 | Phillips |
| M28 | 2.7 | 10g | 28 | Button | $35.66 | No | 1067 | Phillips |
| M36 | 3.2 | 12g | 36 | Pan | $41.32 | No | 434 | Both |
| M50 | 4.5 | 15g | 50 | Pan | $44.72 | No | 740 | Flat |

# DISCRETE AND CONTINUOUS ATTRIBUTES

▸ Discrete

▸ Has only a finite or countably infinite set of values

▸ Examples: zip codes, set of words in a collection of documents

▸ Often represented as integer variables

▸ Continuous

▸ Has real numbers as attribute values

▸ Examples: temperature, height

▸ Continuous attributes are typically represented as floating-point variables

# TABULAR DATA

▸ Collection of records, each of which consists of a fixed set of attributes

| Name | Thread pitch (mm) | Minor diameter tolerance | Nominal diameter (mm) | Head shape | Price for 50 screws | Available at factory outlet? | Number in stock | Flat or Phillips head? |
|------|------|------|------|------|------|------|------|------|
| M4 | 0.7 | 4g | 4 | Pan | $10.08 | Yes | 276 | Flat |
| M5 | 0.8 | 4g | 5 | Round | $13.89 | Yes | 183 | Both |
| M6 | 1 | 5g | 6 | Button | $10.42 | Yes | 1043 | Flat |
| M8 | 1.25 | 5g | 8 | Pan | $11.98 | No | 298 | Phillips |
| M10 | 1.5 | 6g | 10 | Round | $16.74 | Yes | 488 | Phillips |
| M12 | 1.75 | 7g | 12 | Pan | $18.26 | No | 998 | Flat |
| M14 | 2 | 7g | 14 | Round | $21.19 | No | 235 | Phillips |
| M16 | 2 | 8g | 16 | Button | $23.57 | Yes | 292 | Both |
| M18 | 2.1 | 8g | 18 | Button | $25.87 | No | 664 | Both |
| M20 | 2.4 | 8g | 20 | Pan | $29.09 | Yes | 486 | Both |
| M24 | 2.55 | 9g | 24 | Round | $33.01 | Yes | 982 | Phillips |
| M28 | 2.7 | 10g | 28 | Button | $35.66 | No | 1067 | Phillips |
| M36 | 3.2 | 12g | 36 | Pan | $41.32 | No | 434 | Both |
| M50 | 4.5 | 15g | 50 | Pan | $44.72 | No | 740 | Flat |

# DOCUMENT DATA

▸ Each document is represented as a **term** vector, where each attribute records the number of times the term occurs in the document

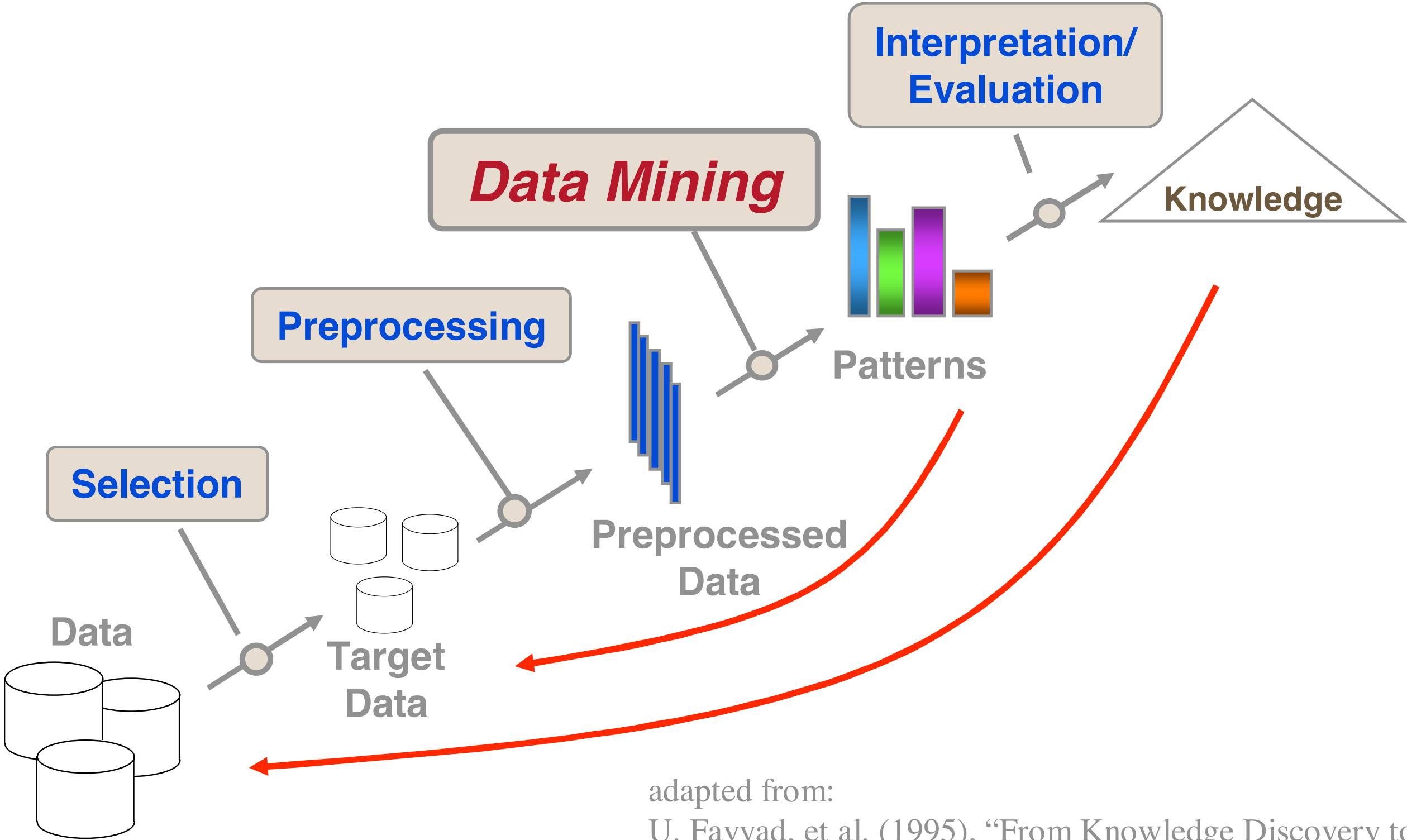| Terms | M1 | M2 | M3 | M4 | M5 | M6 | M7 | M8 | M9 | M10 | M11 | M12 | M13 | M14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| abnormalities | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| age | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| behavior | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| blood | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| close | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| culture | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| depressed | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| discharge | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| disease | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| fast | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 |
| generation | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| oestrogen | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| patients | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| pressure | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 |
| rats | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| respect | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| rise | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| study | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |

# TRANSACTION DATA

▸ Each record corresponds to a transaction that involves a set of items

▸ E.g., in a grocery store purchase, the set of products purchased by a customer constitute a transaction, while the individual products that were purchased are the items

| Customer ID | Transaction ID | Items Bought |
|---|---|---|
| 1 | 0001 | {a,d,e} |
| 1 | 0024 | {a,b,c,e} |
| 2 | 0012 | {a,b,d,e} |
| 2 | 0031 | {a,c,d,e} |
| 3 | 0015 | {b,c,e} |
| 3 | 0022 | {b,d,e} |
| 4 | 0029 | {c,d} |
| 4 | 0040 | {a,b,c} |
| 5 | 0033 | {a,d,e} |
| 5 | 0038 | {a,b,e} |

Table 6.22. Example of market basket transactions.

# ELEMENTS OF DATA MINING ALGORITHMS

# DATA MINING PROCESS



adapted from:
U. Fayyad, et al. (1995), "From Knowledge Discovery to Data
Mining: An Overview," Advances in Knowledge Discovery and
Data Mining, U. Fayyad et al. (Eds.), AAAI/MIT Press

# OVERVIEW

▸ Task specification

▸ Knowledge representation

▸ Learning technique

   ▸ Search + scoring

▸ Prediction and/or interpretation

# OVERVIEW

▸ **Task specification**

▸ Knowledge representation

▸ Learning technique

   ▸ Search + scoring

▸ Prediction and/or interpretation

# TASK SPECIFICATION

▸ Objective of the person who is analyzing the data

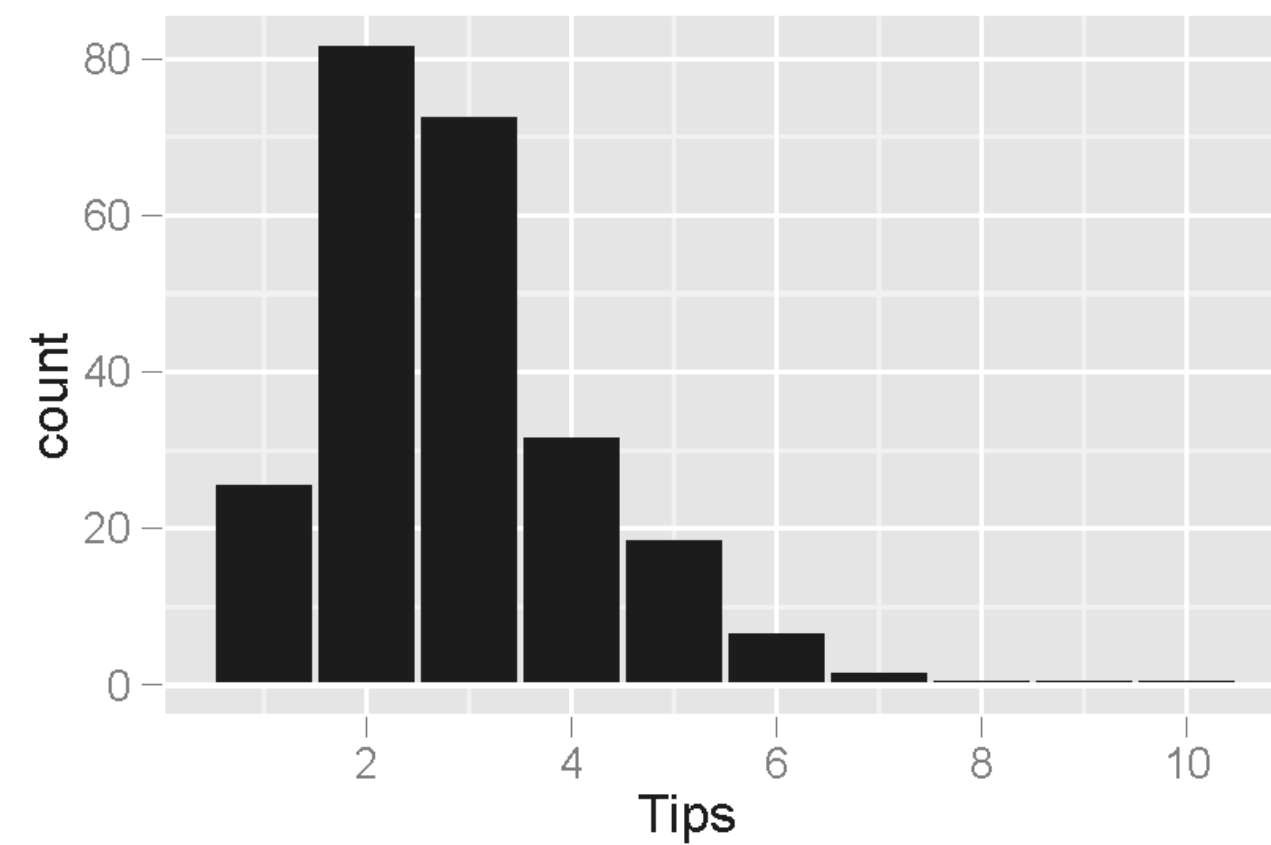▸ Description of the characteristics of the analysis and desired result

# EXPLORATORY DATA ANALYSIS

▸ Goal

  ▸ Interact with data without clear objective

  ▸ Summarize the main characteristics of the data

▸ Techniques

  ▸ Mostly visualization

# EXPLORATORY DATA ANALYSIS EXAMPLE

▸ What influences the amount of tip that a dining party will give to the waiter?
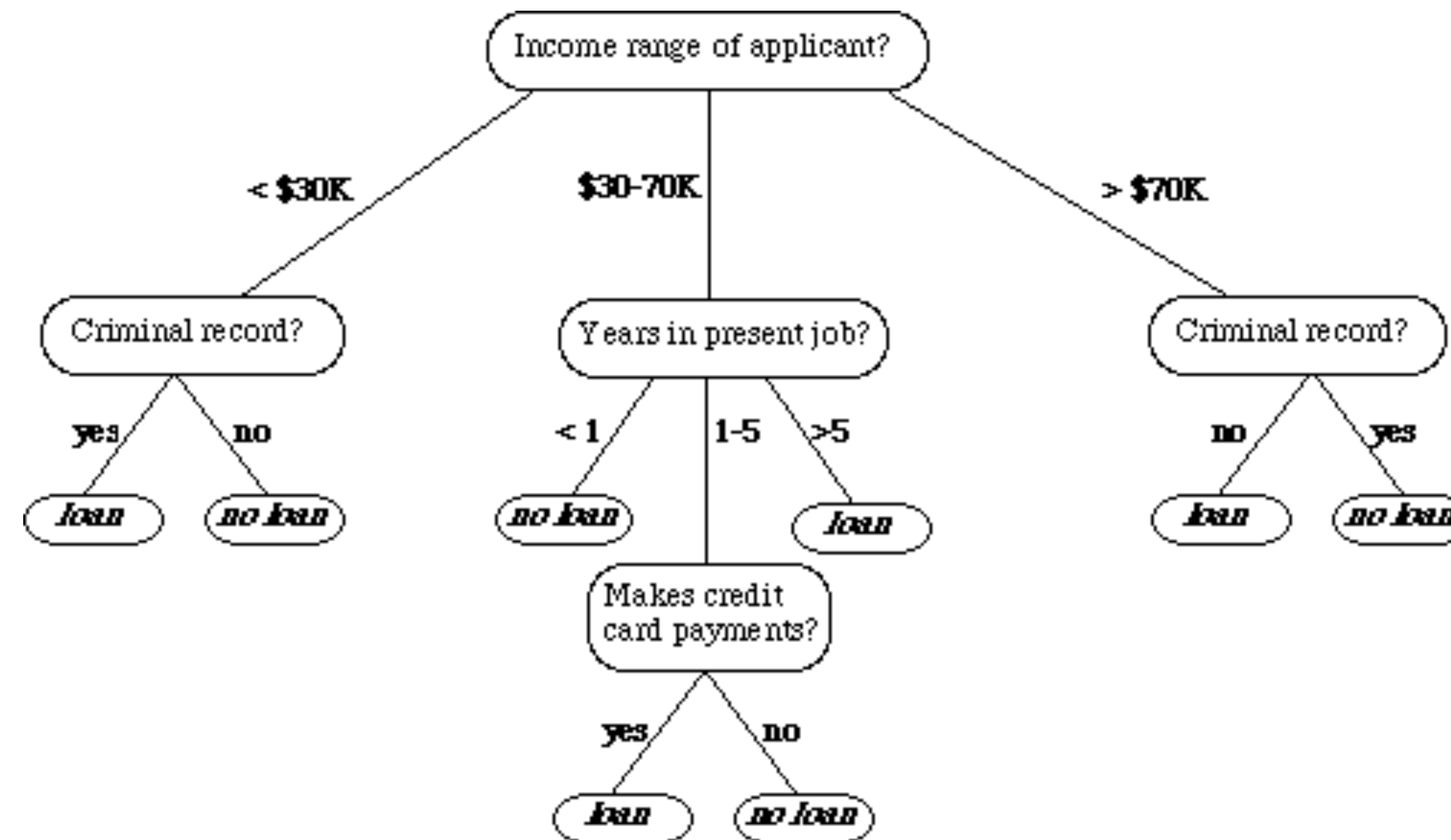


Cook, D. And Dwayne, D. F. Interactive and Dynamic Graphics for Data Analysis: With R and GGobi

# PREDICTIVE MODELING

▸ Goal

  ▸ Learn model to predict the unknown value of a variable of interest given observed attribute values

▸ Techniques

  ▸ Classification, regression

Also known as: **supervised** learning

# PREDICTIVE MODELING EXAMPLE

▸ Zestimate: House sales price prediction!

▸ Predicting loan repayment (and thus decide whether to provide a loan)
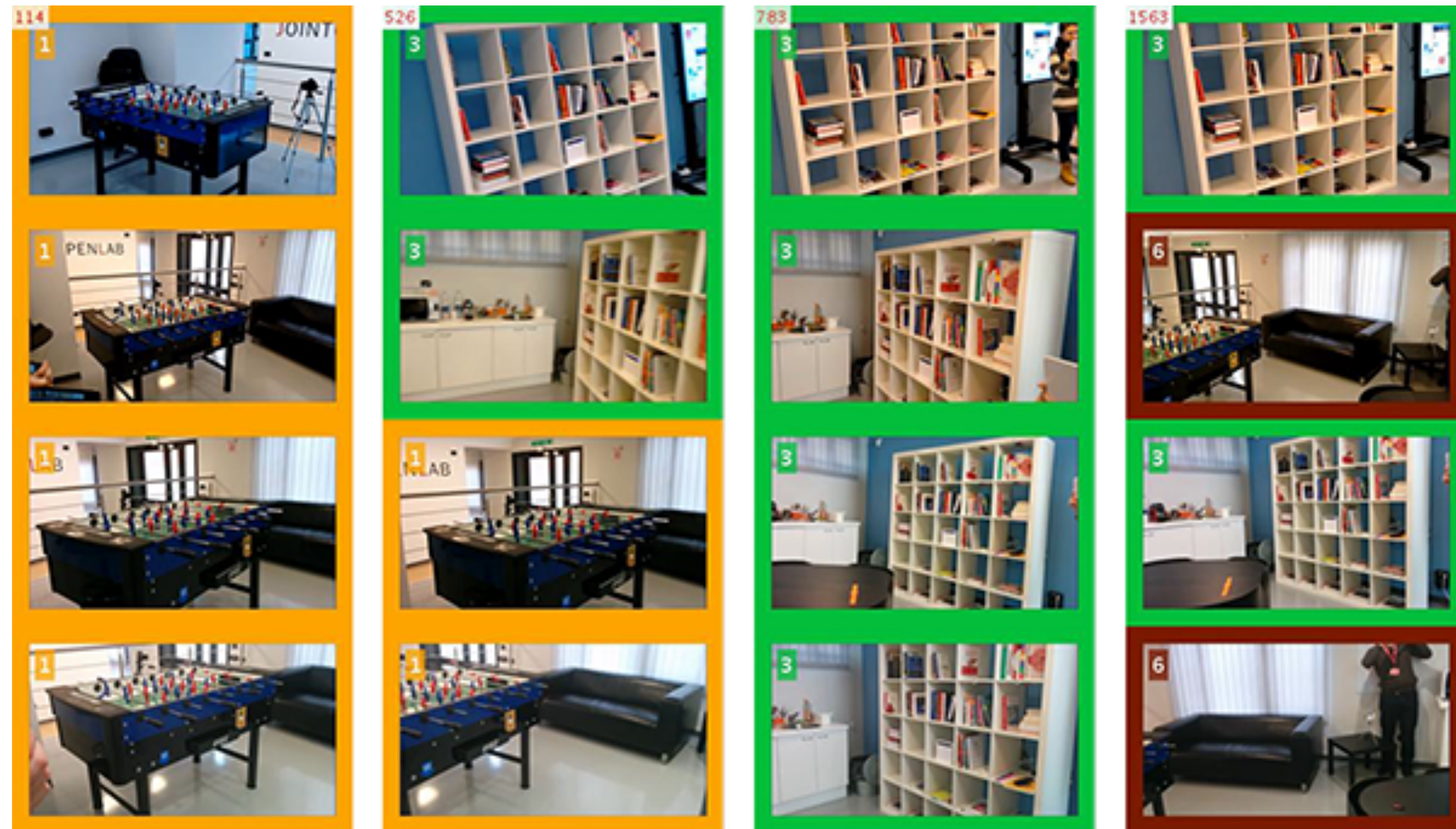
# DESCRIPTIVE MODELING

▸ Goal

  ▸ Summarize the data or the underlying generative process

▸ Techniques

  ▸ Density estimation, cluster analysis and segmentation, probabilistic graphical model

Also known as: **unsupervised** learning

# DESCRIPTIVE MODELING EXAMPLE

▸ Video/scene clustering



Milotta et al. RECfusion: Automatic Scene Clustering and Tracking in Videos from Multiple Sources
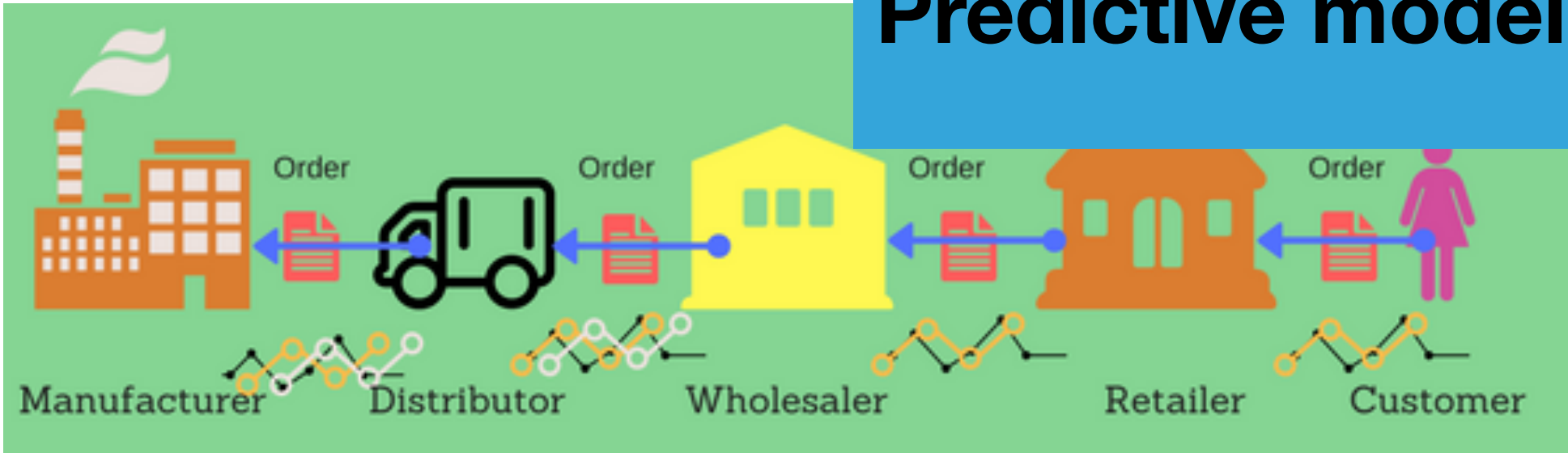
# PATTERN DISCOVERY

▸ Goal

  ▸ Detect patterns and rules that describe subsets of examples

▸ Techniques

  ▸ Association rules, anomaly detection, etc.

**Model**: global summary of a data set
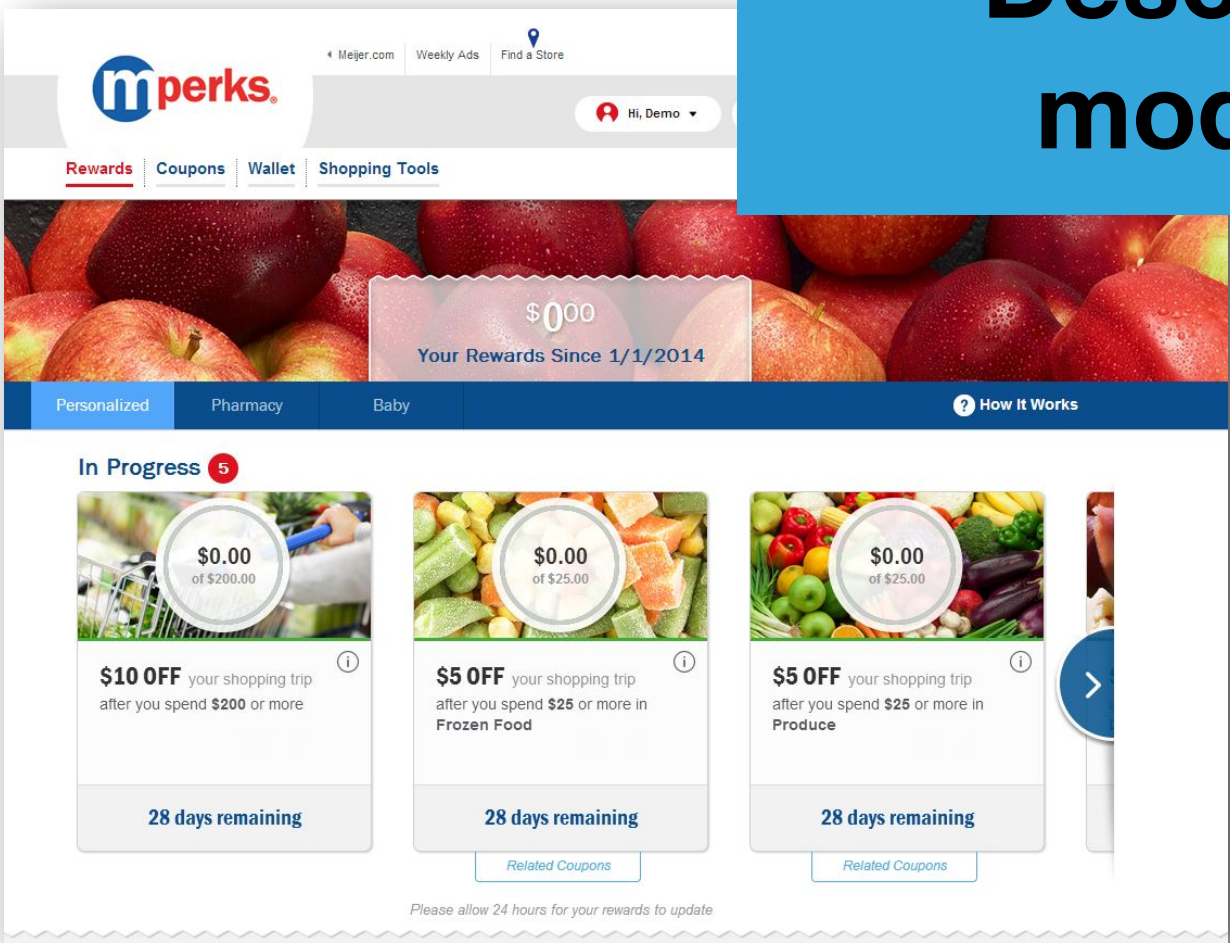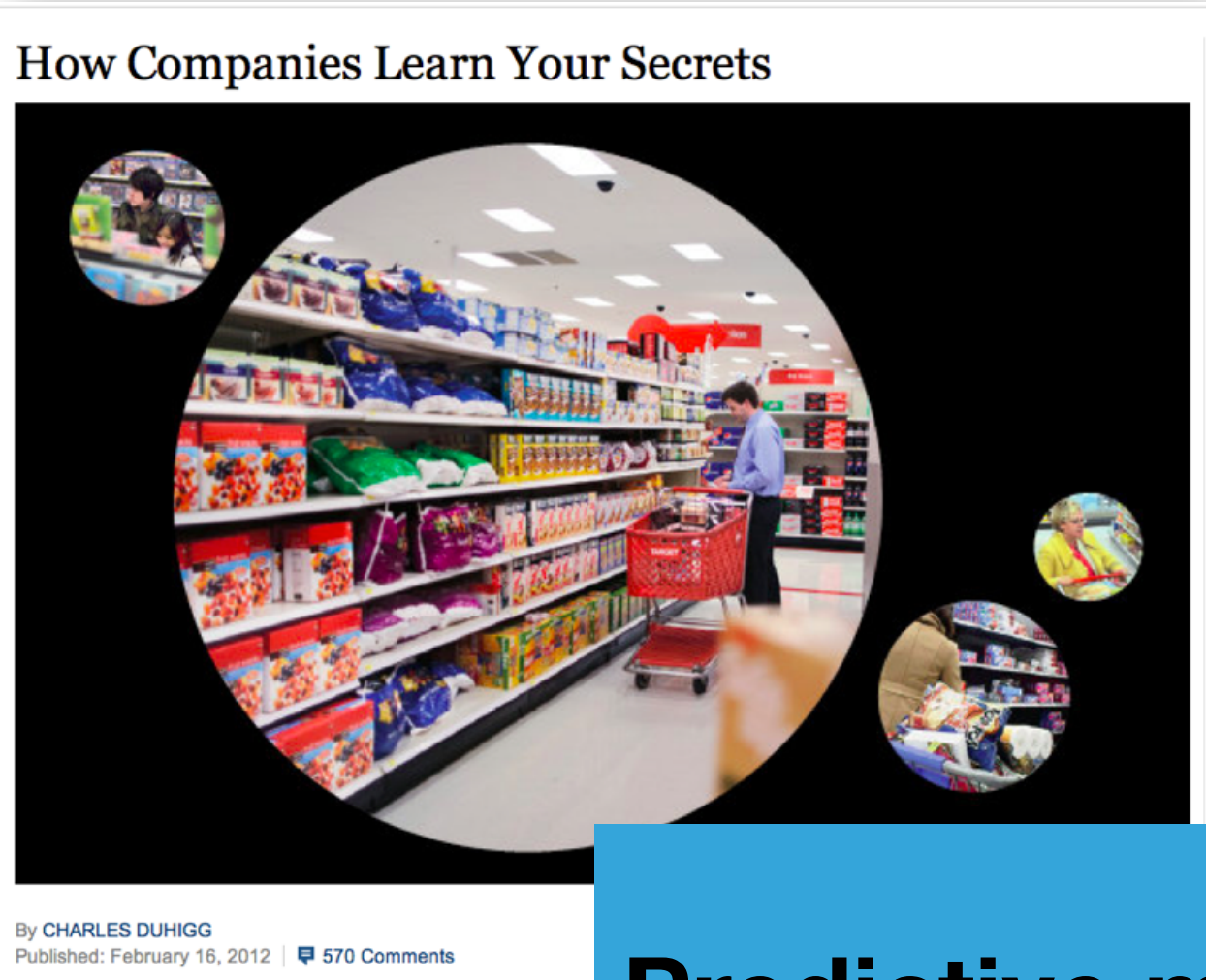**Pattern**: local to a subset of the data

# WHAT TASKS ARE THEY?

**Predictive modeling**

Sales and inventory forecast

**Descriptive modeling**

Customer segmentation

**Predictive modeling**

Pregnant custo

**Pattern Discovery**

Beer & Dia

# OVERVIEW

▸ Task specification

▸ **Knowledge representation**

▸ Learning technique

   ▸ Search + scoring

▸ Prediction and/or interpretation

# KNOWLEDGE REPRESENTATION

▸ Underlying structure of the model or patterns that we seek from the data

  ▸ Specifies the models/patterns that could be returned as the results of the data mining algorithm

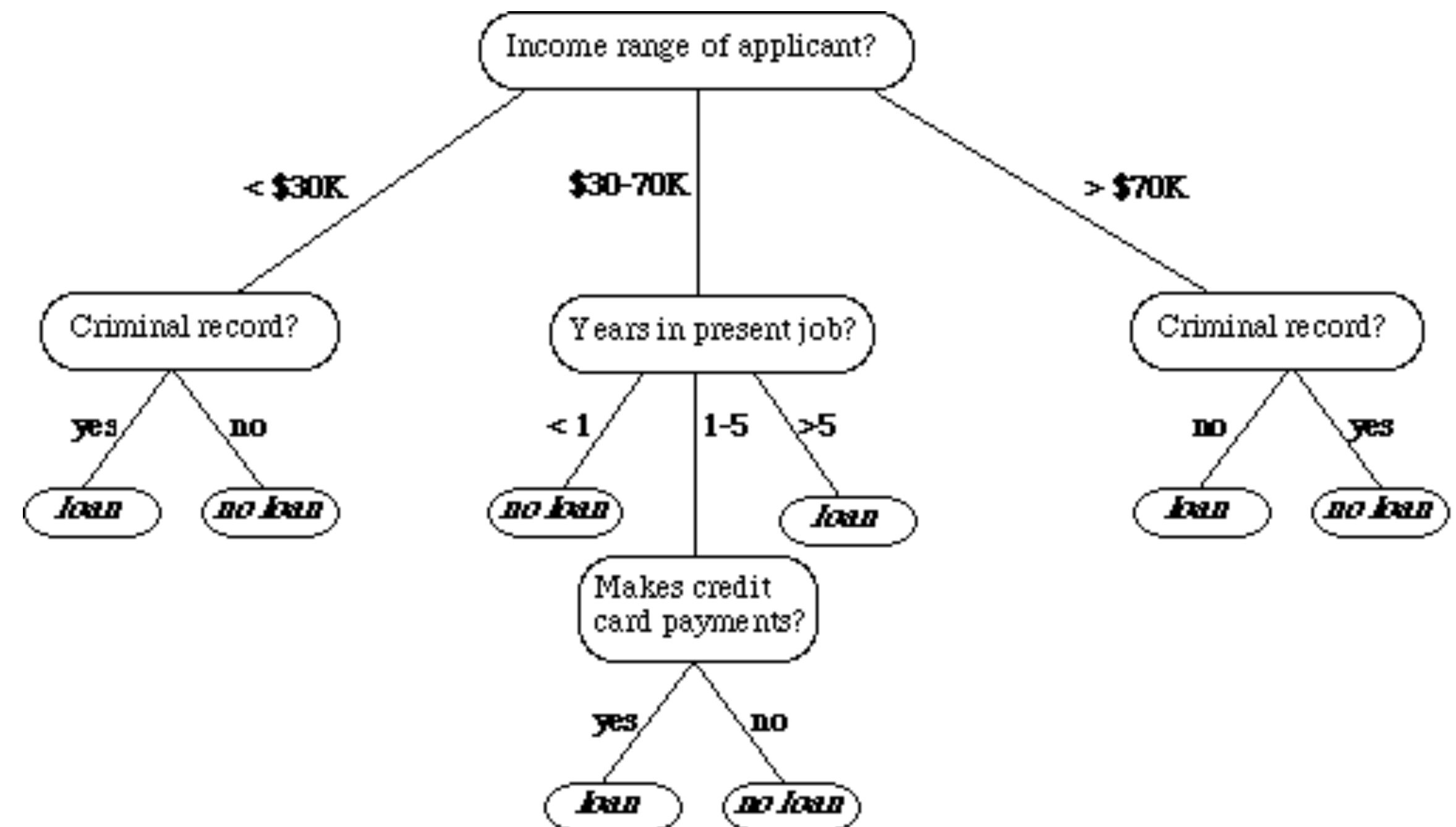  ▸ Defines space of possible models/patterns for algorithm to search over

# KNOWLEDGE REPRESENTATION EXAMPLE: PREDICTIVE MODELING

▸ If-then rule

  ▸ *If* (personal income > $70k) AND (criminal record = 'no'), *then* loan=yes

▸ Decision tree

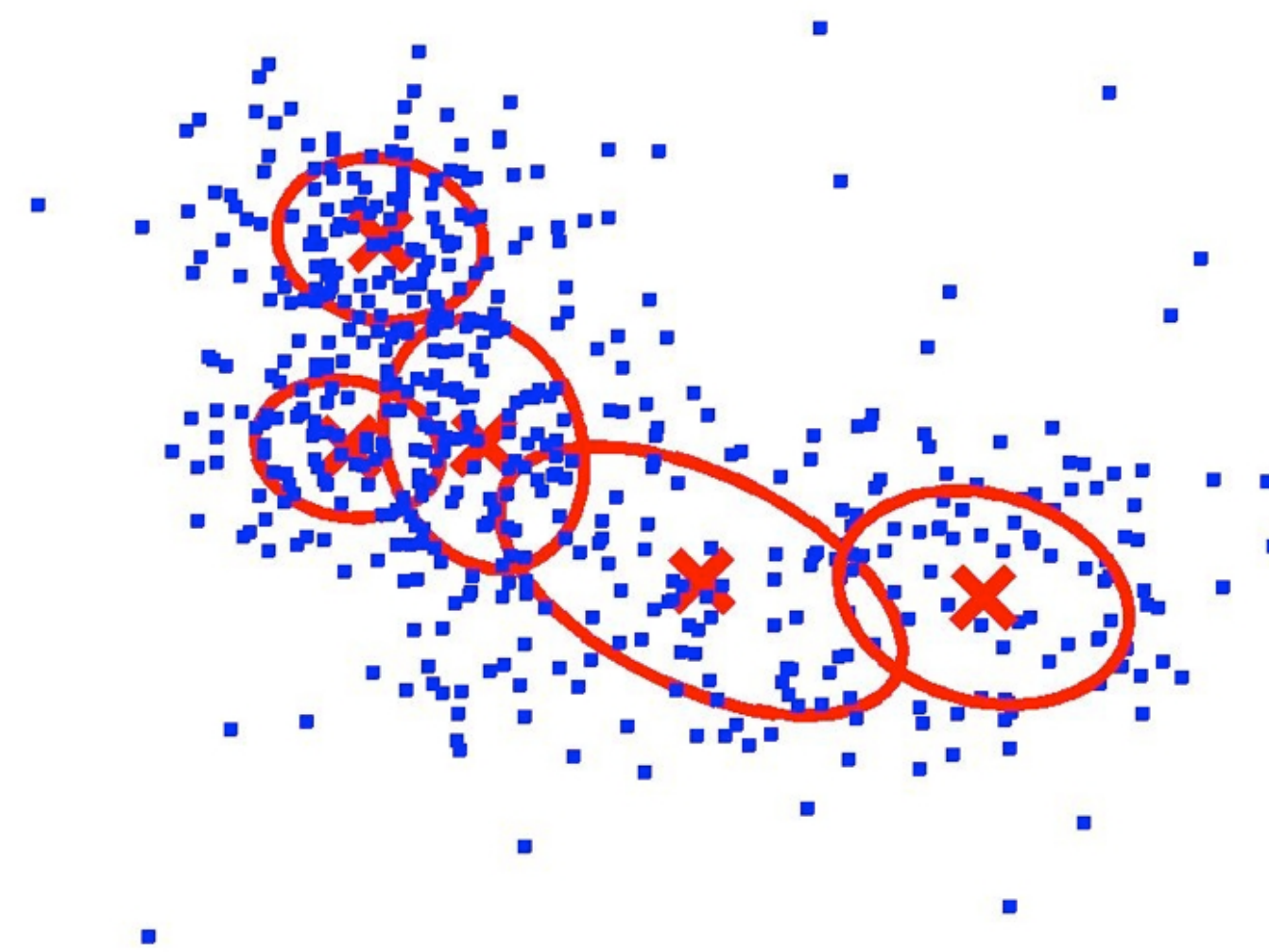  ▸ Each node corresponds to an attribute

  ▸ Each leaf is a class label

# KNOWLEDGE REPRESENTATION EXAMPLE: PREDICTIVE MODELING

▸ Conditional probability distributions (i.e., P($Y$ | $\boldsymbol{X}$))

    ▸ Logistic regression: $log \dfrac{P(Y = 1|\boldsymbol{x})}{1 - P(Y = 1|\boldsymbol{x})} = \beta_0 + \boldsymbol{\beta x}$

    ▸ Model the log-odds as a linear combination of predictors

▸ Linear regression

    ▸ $y = \beta_1 x_1 + \beta_2 x_2 \ldots + \beta_0$

    ▸ $y$ is the response variable, $\boldsymbol{x}$ is the predictor variable

# KNOWLEDGE REPRESENTATION EXAMPLE: DESCRIPTIVE MODELING

▸ Mixture model: Instances represented as a weighted combination of mixture distributions

$$f(x) = \sum_{k=1}^{K} w_k f_k(x; \theta)$$

**likelihood of x being generated from cluster k**

**probability of observing x**

**likelihood of point belonging to cluster k**

# KNOWLEDGE REPRESENTATION EXAMPLE: PATTERN DISCOVERY

▸ Association rules

  ▸ $I = \{i_1, i_2, \ldots, i_n\}$ is a set of $n$ items

  ▸ $T = \{t_1, t_2, \ldots, t_m\}$ is a set of $m$ transactions

**Example database with 5 transactions and 5 items**

| transaction ID | milk | bread | butter | beer | diapers |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 | 0 |
| 2 | 0 | 0 | 1 | 0 | 0 |
| 3 | 0 | 0 | 0 | 1 | 1 |
| 4 | 1 | 1 | 1 | 0 | 0 |
| 5 | 0 | 1 | 0 | 0 | 0 |

▸ An association rule has the form $X{\rightarrow}Y$, where $X$ and $Y$ are subsets of $I$, which means if items in $X$ appear in a transaction, then items in $Y$ are likely to appear in that transaction

  ▸ E.g., {beer} → {diaper}; {bread} → {milk}

# OVERVIEW

▸ Task specification

▸ Knowledge representation

▸ **Learning technique**
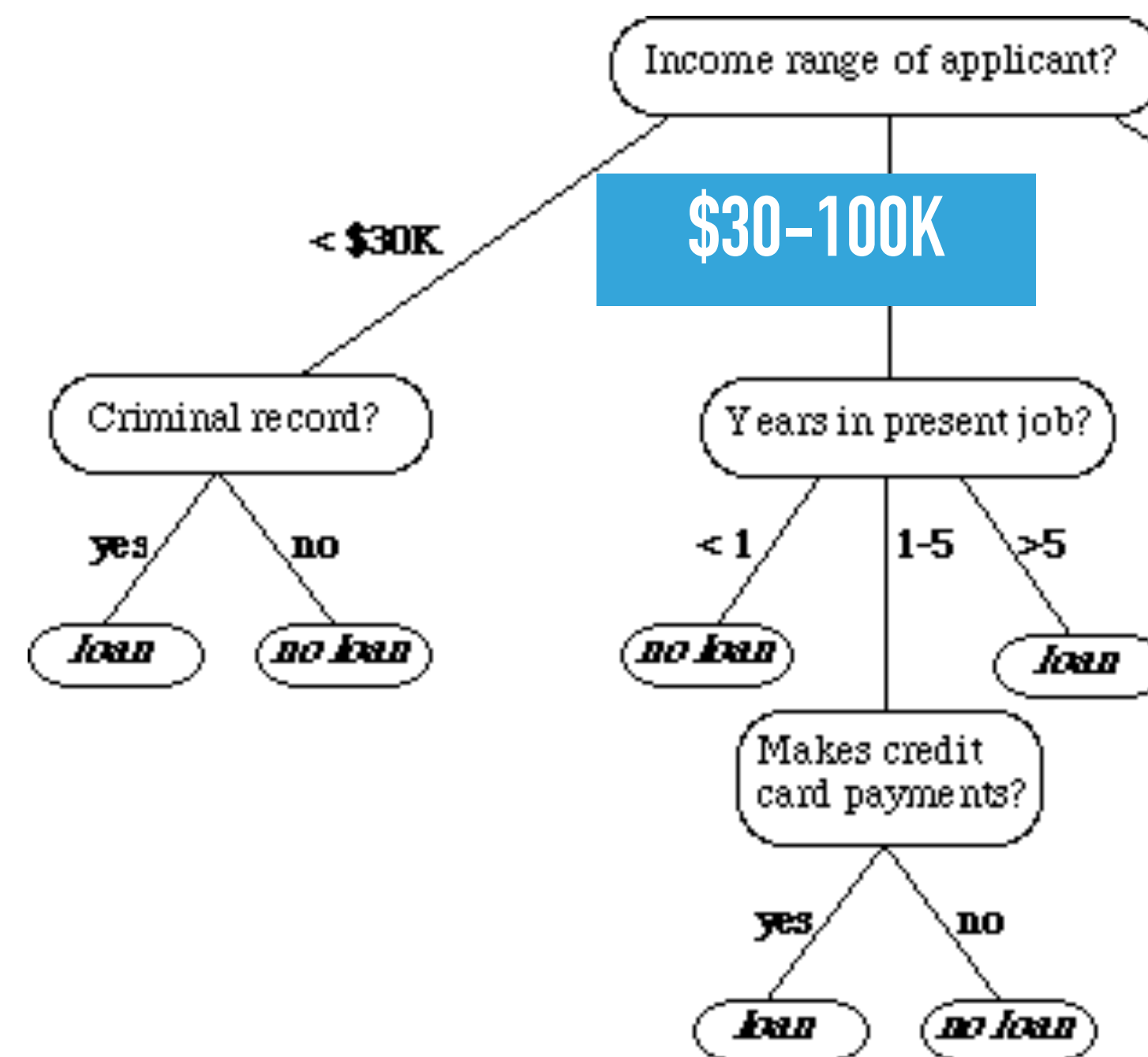
  ▸ Search + scoring

▸ Prediction and/or interpretation

# LEARNING TECHNIQUE

▸ Method to construct model or patterns from data

▸ **Model space**

  ▸ Choice of knowledge representation defines a set of possible models or patterns

▸ **Scoring function**

  ▸ Associates a numerical value (score) with each member of the set of models/patterns

▸ **Search technique**

  ▸ Defines a method for generating members of the set of models/patterns, determining their score, and identifying the ones with the "best" score

# MODEL SPACE

▸ Defined by the choice of knowledge representation

▸ Decision tree:

Income range of applicant?

< $30K     $30-100K     >$100K     **What values to split on?**

Criminal record?     Years in present job?     **Loan**

yes   no     < 1   1-5   >5

loan   no loan     no loan     loan

Makes credit card payments?     **What attributes to include?**

yes   no

loan   no loan

# MODEL PARAMETERS AND STRUCTURE

‣ Models have both **parameters** and **structure**

‣ **Parameters:**

  ‣ Feature values in classification tree

  ‣ Coefficients in regression model

  ‣ Probability estimates in graphical model

‣ **Structure:**

  ‣ Nodes in classification tree

  ‣ Variables in regression model

  ‣ Edges in graphical model

# SCORING FUNCTION

▸ A numeric score assigned to each possible model in a search space, **given a reference/input dataset**

  ▸ Used to judge the quality of a particular model for the domain

▸ Score function are **statistics**—estimates of a population parameter based on a sample of data

▸ Examples:

  ▸ Misclassification

  ▸ Squared error

  ▸ Likelihood

# PARAMETER ESTIMATION VS. STRUCTURE LEARNING

▸ **Parameters:**

   ▸ Feature values in classification tree

   ▸ Coefficients in regression model

   ▸ Probability estimates in graphical model

*Search*: Convex/smooth optimization techniques

▸ **Structure:**

   ▸ Nodes in classification tree

   ▸ Variables in regression model

   ▸ Edges in graphical model

*Search*: Heuristic approaches for combinatorial optimization