# Introduction to Big Data, Predictive Analytics, and Data Science

# Big Data and Data Science Everywhere



Web search and online ads

Insurance

Telcos

Online Education

Online Retail

Social Networks

Entertainment

Healthcare

# Big Data and Data Science Everywhere

AND MANY OTHER PLACES....

# Online Shopping

# Social Networks

# Online Entertainment

# Web Search

# Brainstorming

- What are some other applications?

# Connecting the Dots

- The underlying magic behind what we saw is 'big data' and 'predictive analytics'

# Big Data Pipeline

**Stage: Data influx**
- **Output:** Data stream

**Stage: Collection**
- **Output:** Target data

**Stage: Preprocessing**
- **Output:** Preprocessed data

**Stage: Transformation**
- **Output:** Transformed data

**Stage: Data Mining**
- **Output:** Patterns

**Stage: Interpretation and Evaluation**
- **Output:** Knowledge discovery and actionable insights

Big Data – Technology, Platforms & Products

# Data Mining Tasks

- ## Descriptive Methods:
  - Find human-interpretable patterns that describe the data
  - Techniques: Clustering, Association Analysis, x-point summaries

- ## Predictive Methods:
  - Use available data to build models that can predict the outcome of future data
  - Techniques: Classification, Regression, Anomaly, and Deviation Detection

- ## Prescriptive Methods:
  - Predict future outcomes and suggest actions that may prevent or mitigate the impact of the predicted outcomes
  - Techniques: Various optimization techniques

# Traffic Management



## Descriptive [Informing Role]:

- Traffic jam has happened already.
- [Implicit: Do something about it.]

# Traffic Management



## Predictive [Informing and Warning Role]:

- Traffic jam is about to happen in the next 30 minutes.
- [Implicit: Do something before it happens.]

# Traffic Management



## Prescriptive [Informing, Warning, and Advisory Role]:

Take action so traffic jam does not happen

OR

Traffic jam is about to happen in the next 30 minutes and you could possibly take the following courses of action:

- Route traffic to service road near I-5

- Block more traffic from entering the WA-520 bridge

# Online Travel



**Descriptive Analytics:** Historical price trend and variation
**Predictive Analytics:** Price may rise in next 7 days
**Prescriptive Analytics:** *Advice: Buy Confidence:* 85%

**KAYAK**    Flights    Hotels    Cars    Deals    Vacations

Seattle (SEA)    ▶    San Francisco (SFO)    04/17/2013    ▶    04/24/2013

Hide toolbox ▼

Price alert        Fare charts
Airline fees       Add baggage
Airline Matrix     +/- 3 days

**Price Trend**

Advice: **Buy**    Confidence: **85%**
Prices may rise within 7 days ⓘ

**308** of **471 flights**    show all

Sort by

**$208**
& up

Select

**Virgin America**

$208 nonstop

Fly with WiFi, on-demand food, live TV, movies, music, and more

Virgin America

# Data Mining and Predictive Analytics

In the next few slides, we will take a look at some of the most common data mining tasks.

# Classification: A Simple Example

|  |  | categorical | categorical | continuous | class |
|--|--|--|--|--|--|

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Single | 75K | ? |
| Yes | Married | 50K | ? |
| No | Married | 150K | ? |
| Yes | Divorced | 90K | ? |
| No | Single | 40K | ? |
| No | Married | 80K | ? |

**Test Set**

**Training Set** → **Learn Classifier** → **Model**

# Classification

- Given a collection of records (training set)
  - Each record contains a set of *attributes*; one of the attributes is the *class label*.
- Find a model for class attribute as a function of the values of other attributes.
- Goal: previously unseen records should be assigned a class as accurately as possible.

# Classification: More Examples

- ## Direct Marketing
  - Goal: reduce cost of mailing by targeting a set of consumers likely to buy a new cell-phone product
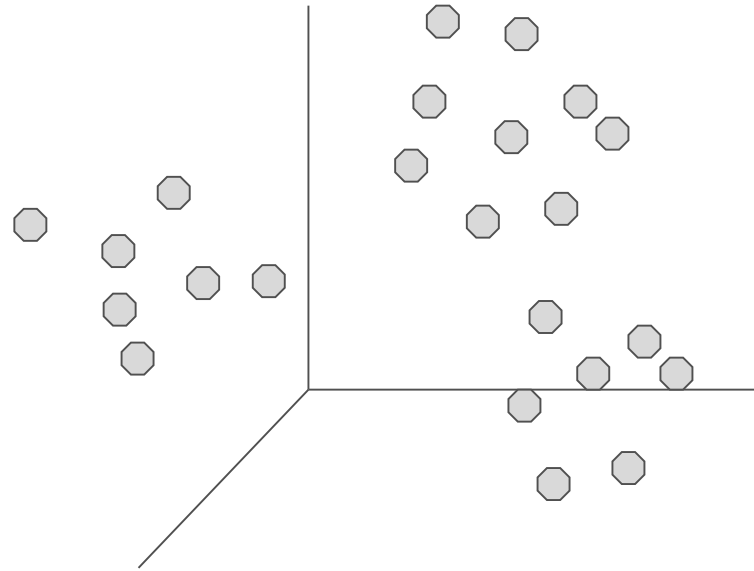
- ## Fraud Detection
  - Goal: predict fraudulent cases in credit card transactions
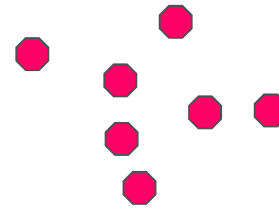
- ## Customer Attrition/Churn
  - Goal: predict whether a customer is likely to be lost to a competitor
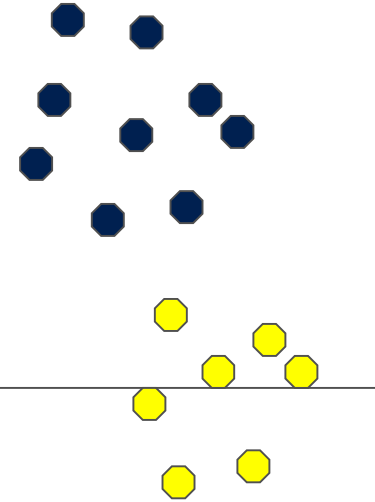
# Clustering: An Illustration

Clustering in 3-D space using Euclidean distance

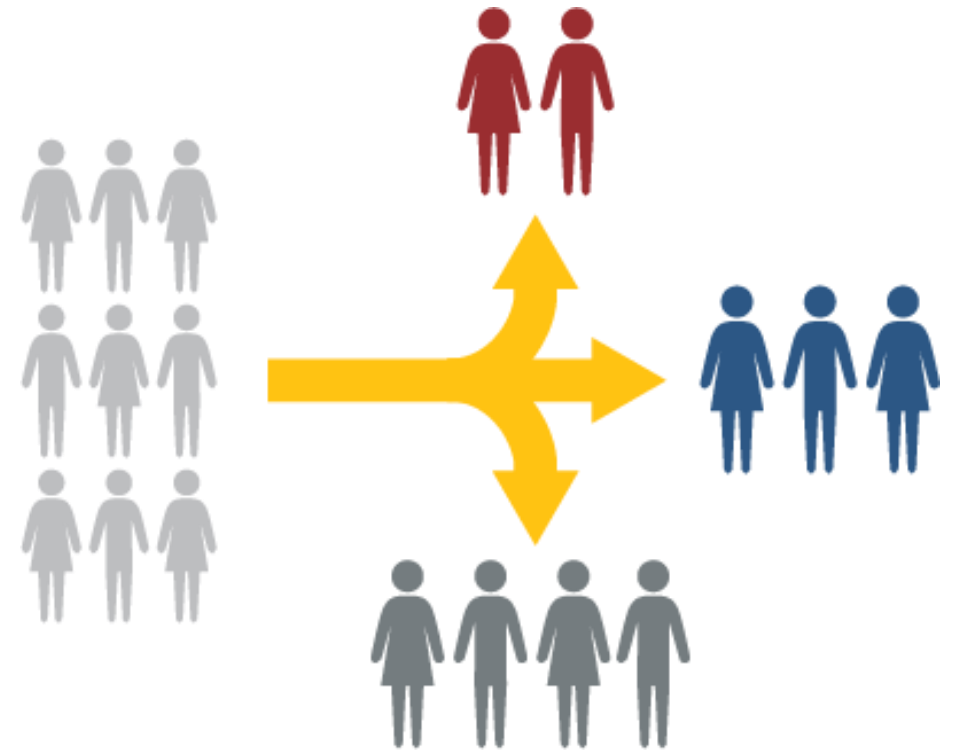Intra-cluster distances
are minimized

Inter-cluster distances
are maximized

# Clustering: Examples

- Subdivide the market into distinct subsets of customers where any subset may conceivably be selected as a segment to be reached with a particular offer
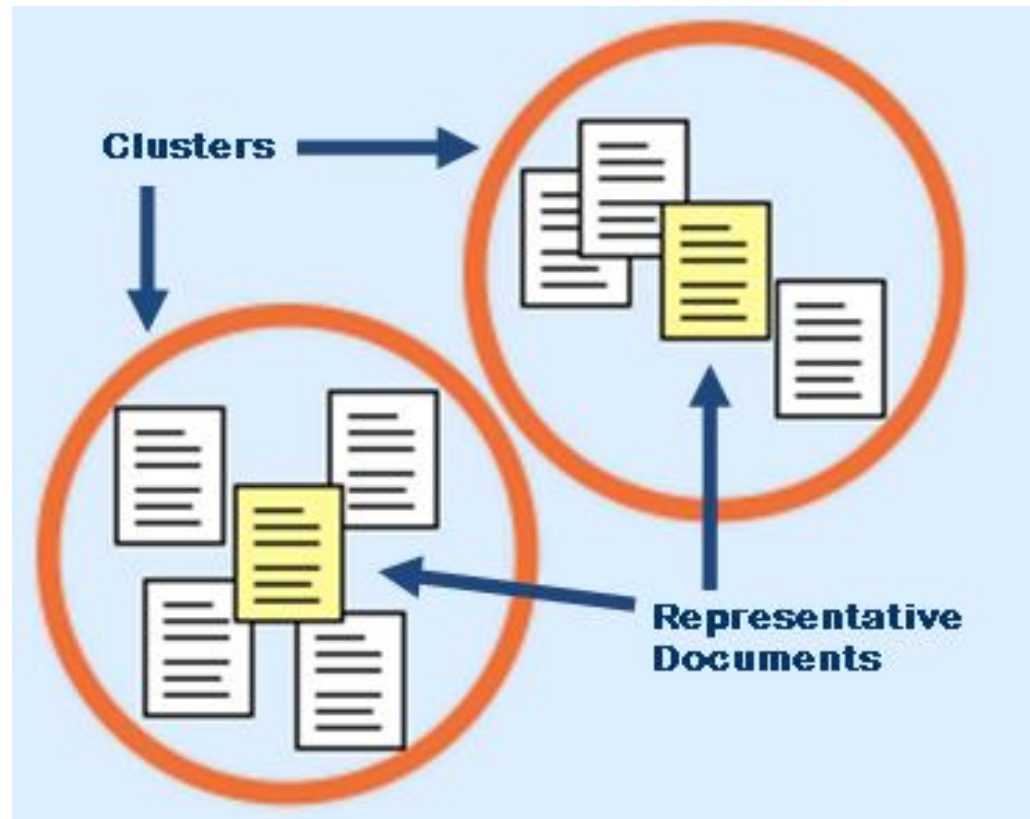
# Clustering

- Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters such that:
    - Data points within a cluster have more similarities with one another
    - Data points in different clusters have less similarities with one another
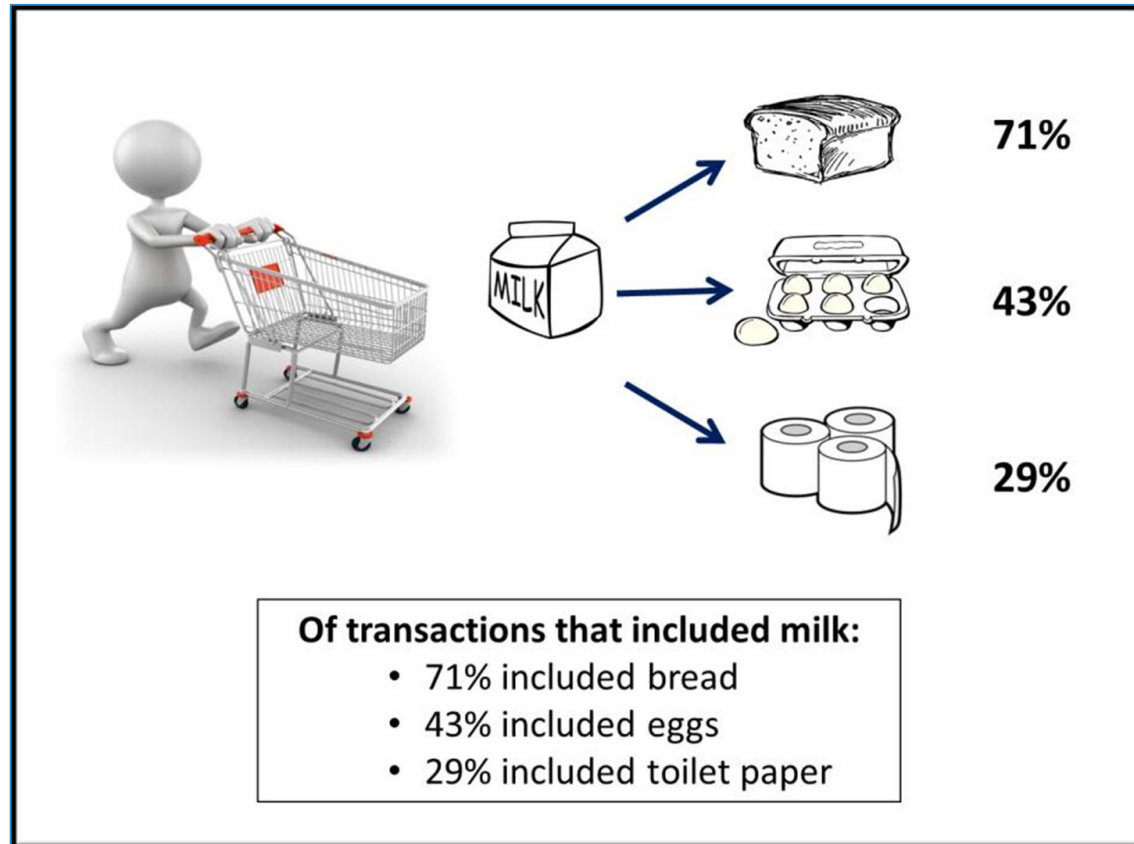
# Clustering: Similarity Measures

- Similarity Measures:
  - Euclidean Distance if attributes are continuous
  - Other problem-specific measures
  - Example: If a particular word occurs in two documents or not

# Clustering: Examples

- To find groups of documents that are similar to each other based on the important terms appearing in them

# Association Analysis



Of transactions that included milk:
- 71% included bread
- 43% included eggs
- 29% included toilet paper

Your behavior is being predicted, not by studying you, but by studying others.

# Association Rule Discovery

- Given a set of records each of which contain some number of items from a given collection:
  - Produce dependency rules which will predict the occurrence of an item based on the occurrences of other items

| TID | Items |
|-----|-------|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

Rules Discovered:
   {Milk} --> {Coke}
   {Diaper, Milk} --> {Beer}

# Association Analysis: Supermarket Shelf Management

- Goal: To identify items that are bought together by a sufficient amount of customers

- Place the items close to each other on supermarket shelves

# Association analysis examples

- ## Marketing and sales promotion:
  - Users who buy item A usually also buy item B
  - If users bought item A, suggest item B or even offer discount on item B

- ## Inventory management:
  - Goal: A consumer appliance repair company wants to anticipate the nature of repairs on its consumer products and keep the service vehicles equipped with the right parts to reduce the number of visits to consumer households

# Regression Example: Predict Housing Prices

# Regression

- Predict a value of a given continuous valued variable based on the values of other variables, assuming a linear or nonlinear model of dependency

# Regression: Ad Clicks

Predict the probability of whether or not an ad will be clicked

# Deviation/Anomaly Detection

- Detect significant deviations from normal behavior

- Applications:
  - Credit Card Fraud Detection
  - Network Intrusion Detection
  - Bot detection in web traffic

Estimate: 50-61% of web traffic is bots

# Challenges in Data Mining

Scalability

Dimensionality

Complex and heterogeneous data

Data quality

Data ownership and distribution

Privacy

Reaction time

Many other domain specific issues

# 5 V$_S$ Of Big Data

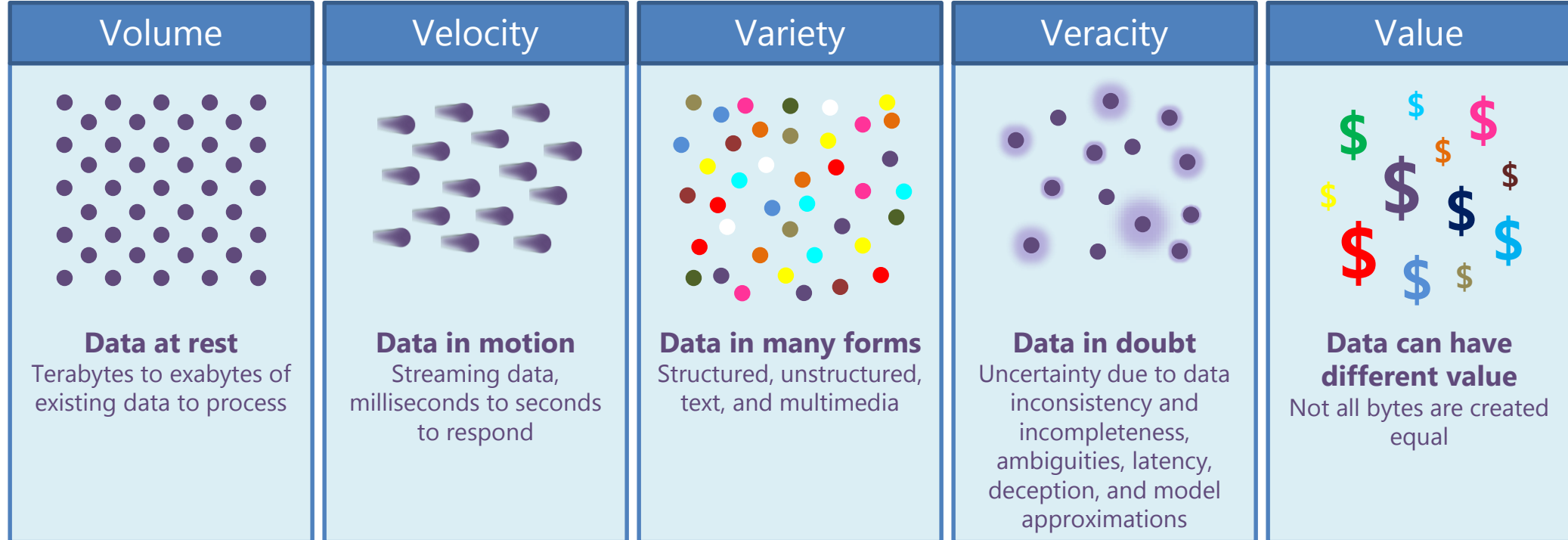| Volume | Velocity | Variety | Veracity | Value |
|---|---|---|---|---|
| **Data at rest**<br>Terabytes to exabytes of existing data to process | **Data in motion**<br>Streaming data, milliseconds to seconds to respond | **Data in many forms**<br>Structured, unstructured, text, and multimedia | **Data in doubt**<br>Uncertainty due to data inconsistency and incompleteness, ambiguities, latency, deception, and model approximations | **Data can have different value**<br>Not all bytes are created equal |

Questions?