

# Lecture 1: Introduction to Business Intelligent Analytics

## What are Data Analytics?

- **Analytic** is the use of:
  - Data,
  - Information technology,
  - Statistical analysis,
  - Quantitative methods,
  - Mathematical or computer-based models.
- To help managers gain improved insight about their business operations and make better, fact-based decisions.
- **Business Analytics (BA) is a subset of Data Analytics.**

## Business Analytics applications:

- Management of customer relationships.
- Financial and marketing activities.
- Supply chain management.
- Human resource planning.
- Pricing decisions.
- Sport team game strategies.

## Importance of Business Analytics:

- There is a strong relationship of BA with:
  - Profitability of businesses.
  - Revenue of businesses.
  - Shareholder return.
- BA enhances understanding of data.
- BA is vital for businesses to remain competitive.
- BA enables creation of informative reports.

## Scope of Business Analytics

- **Descriptive analytics:**
  - Uses data to understand past and present.
- **Predictive analytics:**
  - Analyzes past performance.
- **Prescriptive analytics:**
  - Uses optimization techniques.

## Scope of Business Analytics - Retail Markdown Decisions

- Most department stores clear seasonal inventory by reducing prices.
- The question is:
  - When to reduce the price and how much?

- **Descriptive analytics**: examine historical data for similar products (prices, unit sold, advertising, etc.)
- **Predictive analytics**: predict sales based on price.
- **Prescriptive analytics**: find the best sets of pricing and advertising to maximize sales revenues.

## Data for Business Analytics

- DATA:
  - Collected facts and figures.
- DATABASE:
  - Collection of computer files containing data.
- INFORMATION:
  - Comes from analyzing data.

## Data for Business Analytics

- **Metrics** are used to quantify performance.
- **Measures** are numerical values of metrics.
- **Discrete** metrics involve counting.
  - On time or not on time.
  - Number or proportion of on time deliveries.
- **Continuous** metrics are measured on a continuum.
  - Delivery time.
  - Package weight.
  - Purchase price.

## Data for Business Analytics

- A Sales Transaction Database File:

The diagram illustrates a Sales Transaction Database File. At the bottom, a horizontal bracket spans across the columns, labeled "Fields or Attributes". Above this, a vertical bracket on the left spans across the rows, labeled "Entities". To the right of the table, a large curly brace encompasses all rows and columns, labeled "Records".

A	B	C	D	E	F	G	H	
1	Sales Transactions: July 14							
2								
3	Cust ID	Region	Payment	Transaction Code	Source	Amount	Product	Time Of Day
4	10001	East	Paypal	93816545	Web	\$20.19	DVD	22:19
5	10002	West	Credit	74083490	Web	\$17.85	DVD	13:27
6	10003	North	Credit	64942368	Web	\$23.98	DVD	14:27
7	10004	West	Paypal	70560957	Email	\$23.51	Book	15:38
8	10005	South	Credit	35208817	Web	\$15.33	Book	15:21
9	10006	West	Paypal	20978903	Email	\$17.30	DVD	13:11
10	10007	East	Credit	80103311	Web	\$177.72	Book	21:59
11	10008	West	Credit	14132683	Web	\$21.76	Book	4:04
12	10009	West	Paypal	40128225	Web	\$15.92	DVD	19:35
13	10010	South	Paypal	49073721	Web	\$23.39	DVD	13:26

Figure 1.1

## What is Big Data?

- Information from multiple internal and external sources:
  - Transactions.
  - Social media.
  - Enterprise content.

- Sensors.
- Mobile devices.
- Companies leverage data to adapt products and services to:
  - Meet customer needs.
  - Optimize operations.
  - Optimize infrastructure.
  - Find new source of revenue.
  - Can reveal more patterns and anomalies.
- IBM estimates that by 2015 4.4 million jobs will be created globally to support Big Data.
  - 1.9 million of these jobs will be in United States.

## Types of Data

- When collecting or gathering data we collect data from individuals' cases on variables.
- A **variable** is a unit of data collection whose value can vary.
- Variables can be defined into **types** according to the level of mathematical scaling that can be carried out on the data.
- There are four types of data or levels of measurement:

<b>1. Categorical (Nominal)</b>	<b>2. Ordinal</b>
<b>3. Interval</b>	<b>4. Ratio</b>

## Data for Business Analytics

- **Classifying Data Elements in a Purchasing Database**

A	B	C	D	E	F	G	H	I	J
<b>1 Purchase Orders</b>									
<b>2</b>									
3 Supplier	Order No	Item No.	Item Description	Item Cost	Quantity	Cost per order	A/P Terms (Months)	Order Date	Arrival Date
4 Spacetime Technologies	A0111	6489	O-Ring	\$ 3.00	900	\$ 2,700.00	25	10/10/11	10/18/11
5 Steelpin Inc	A0115	5319	Shielded Cable/ft.	\$ 1.10	17,500	\$ 19,250.00	30	08/20/11	08/31/11
6 Steelpin Inc	A0123	4312	Bolt-nut package	\$ 3.75	4,250	\$ 15,937.50	30	08/25/11	09/01/11
7 Steelpin Inc	A0204	5319	Shielded Cable/ft.	\$ 1.10	16,500	\$ 18,150.00	30	09/15/11	10/05/11
8 Steelpin Inc	A0205	5677	Side Panel	\$195.00	120	\$ 23,400.00	30	11/02/11	11/13/11
9 Steelpin Inc	A0207	4312	Bolt-nut package	\$ 3.75	4,200	\$ 15,750.00	30	09/01/11	09/10/11
10 Alum Sheeting	A0223	4224	Bolt-nut package	\$ 3.95	4,500	\$ 17,775.00	30	10/15/11	10/20/11
11 Alum Sheeting	A0433	5417	Control Panel	\$255.00	500	\$ 127,500.00	30	10/20/11	10/27/11
12 Alum Sheeting	A0443	1243	Airframe fasteners	\$ 4.25	10,000	\$ 42,500.00	30	08/08/11	08/14/11
13 Alum Sheeting	A0446	5417	Control Panel	\$255.00	400	\$ 103,530.00	30	09/01/11	09/10/11
14 Spacetime Technologies	A0533	9752	Gasket	\$ 4.05	1,500	\$ 6,075.00	25	09/20/11	09/25/11
15 Spacetime Technologies	A0555	6489	O-Ring	\$ 3.00	1,100	\$ 3,300.00	25	10/05/11	10/10/11

Figure 1.2      Categorical      Categorical      Categorical      Categorical      Ratio      Ratio      Ratio      Ratio      Interval      Interval

## Categorical (Nominal) data

- **Nominal or categorical** data is data that comprises of categories that **cannot** be rank ordered – each category is just different.
- The categories available **cannot be placed in any order** and no judgement can be made about the relative size or distance from one category to another.
  - Categories bear no quantitative relationship to one another.

- Examples:
  - Customer's location (American, Europe, Asia).
  - Employee classification (manager, supervisor, associate).
- What does this mean? No mathematical operations can be performed on the data relative to each other.
- Therefore, nominal data reflect qualitative differences rather than quantitative.

**What is your gender?  
(please tick)**

Male	<input type="checkbox"/>
Female	<input type="checkbox"/>

**Did you enjoy the film? (please tick)**

Yes	<input type="checkbox"/>
No	<input type="checkbox"/>

- Systems for measuring nominal data must ensure that each category is mutually exclusive, and the system of measurement needs to be exhaustive.
- Variables that have only 2 responses, i.e., Yes or No, are known as dichotomies.

## Ordinal data

- Ordinal data is data that comprises of categories that can be ranked ordered.
- Similarly with nominal data, the distance between each category cannot be calculated but the categories can be ranked above or below each other.
  - No fixed units of measurement.
  - Examples:
    - College football rankings.
    - Survey responses (poor, average, good, very good, excellent).
- What does this mean? Can make statical judgements and perform limited math.

## How satisfied are you with the level of service you have received? (please tick)

Very satisfied

Somewhat satisfied

Neutral

Somewhat dissatisfied

Very dissatisfied

### Interval and ratio data

- Both interval and ratio are examples of **scale data**.
- Scale data:
  - o Data is in numeric format (\$50, \$100, \$200).
  - o Data that can be **measured on a continuous scale**.
  - o The **distance** between each can be observed and as a result **measured**.
  - o The data can be **placed in rank order**.

### Interval data

- Ordinal data but with constant differences between observations.
- Ratios are not meaningful.
- Examples:
  - o **Time** – moves along a continuous measure of seconds, minutes, and so on, and is without a zero point of time.
  - o **Temperature** – moves along a continuous measure of degrees and is without a true zero.
  - o **SAT scores**.

### Ratio data

- Ratio data measured on a **continuous scale** and **does have a natural zero point**.
  - o Ratios are meaningful.
  - o Examples:
    - Monthly sales.
    - Delivery times.
    - Weight.
    - Height.

- Age.

## Types of analytics – Decision Models

- Model:
  - An abstraction or representation of a real system, idea, or object.
  - Captures the most important features.
  - Can be:
    - written or verbal description,
    - a visual display,
    - a mathematical formula,
    - a spreadsheet representation.

- **Decision Models:**



- Is a model used to understand, analyze, or facilitate decision making.
- Types of model input:
  - Data.
  - Uncontrollable variables.
  - Decision variables (controllable).
- **Descriptive Decision Models:**
  - Simply tell “what is” and describe relationships.
  - Do not tell managers what to do.

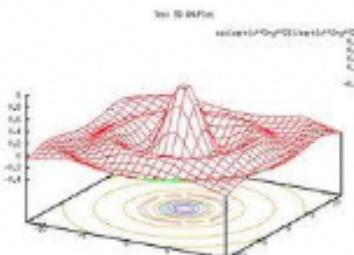
## Descriptive Analytics

- Description analytics, such as reporting/OLAP, dashboards, and data visualization, have been widely used for some time.
- They are the core of traditional BI.

Year 2000		Audit Division		Wales Division	
Line Item		Budget	Actual	Budget	Actual
Cost of Goods Sold		\$6,051,306.49	\$7,112,961.38	\$4,322,514.74	\$4,526,954.71
Marketing Expense		\$750,179.25	\$105,586.17	\$455,040.05	\$462,015.43
Research and Development Expense		\$530,243.39	\$630,014.73	\$129,000.95	\$136,000.13
Selling Expense		\$1,022,301.64	\$1,579,786.18	\$896,887.48	\$927,970.90
taxes		\$314,655.00	\$319,380.18	\$202,636.67	\$200,206.01

Year 2001		Audit Division		Wales Division	
Line Item		Budget	Actual	Budget	Actual
Cost of Goods Sold		\$2,564,566.31	\$2,700,773.16	\$1,726,031.16	\$1,773,448.00
Marketing Expense		\$294,766.23	\$296,686.70	\$167,757.28	\$170,778.55
Research and Development Expense		\$230,719.93	\$193,236.83	\$104,270.95	\$125,725.00
Selling Expense		\$635,427.30	\$611,649.47	\$405,080.93	\$400,161.91
taxes		\$136,886.70	\$122,526.31	\$80,480.78	\$80,671.87

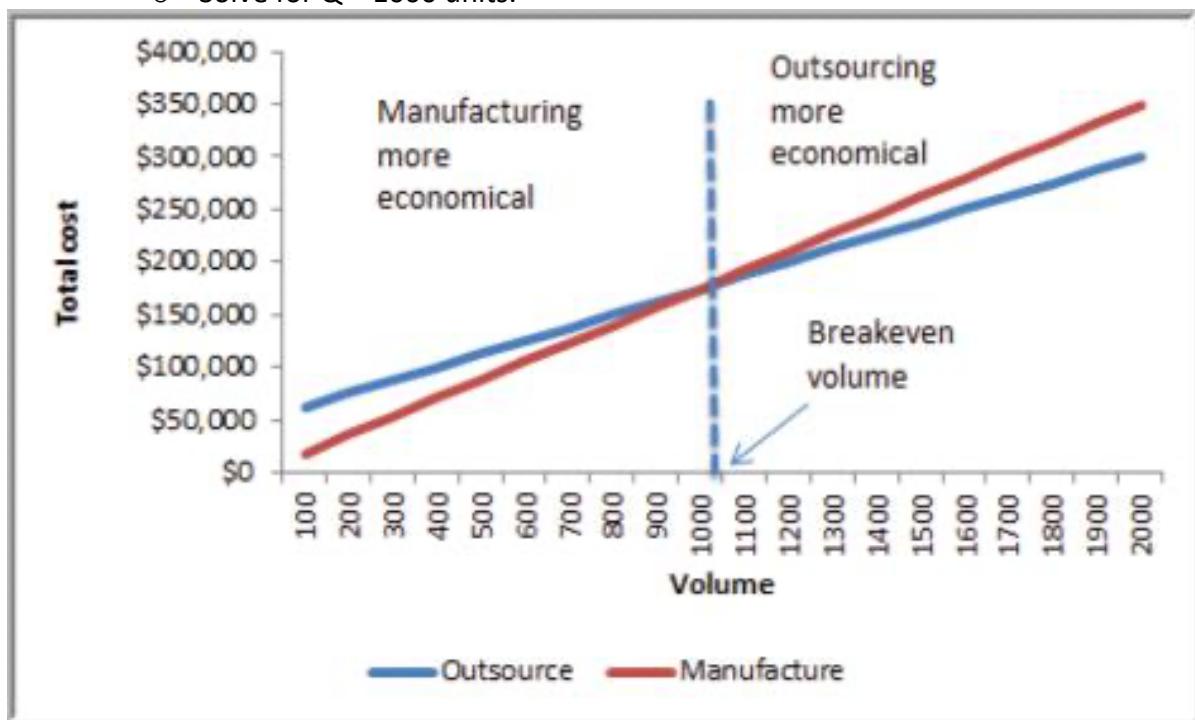


## What has occurred?

Descriptive analytics, such as data visualization, is important in helping users interpret the output from predictive and predictive analytics.

## Decision Models

- A Break-even Decision Model:
  - o  $TC(\text{manufacturing}) = \$50,000 + \$125 * Q$
  - o  $TC(\text{outsourcing}) = \$175 * Q$
- Break-even Point:
  - o Set  $TC(\text{manufacturing}) = TC(\text{outsourcing})$ .
  - o Solve for  $Q = 1000$  units.

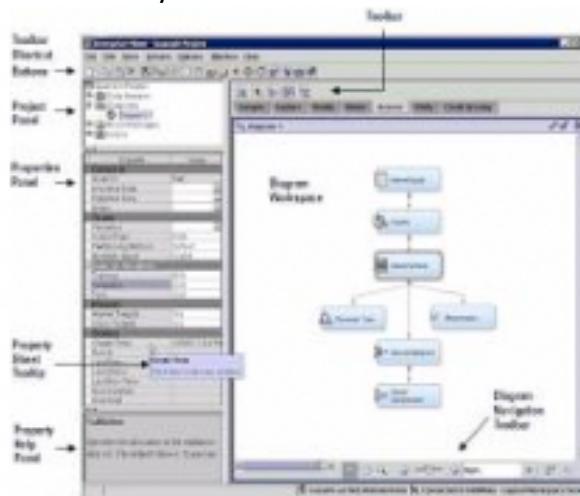
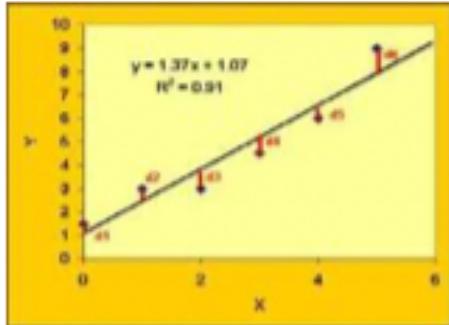


## Predictive Decision Models

- Predictive Decision Models often incorporate uncertainty to help managers analyze risk.
- Aim to predict what will happen in the future.
- Uncertainty is imperfect knowledge of what will happen in the future.
- Risk is associated with the consequences of what happens.

## Predictive Analytics

- Algorithms for predictive analytics, such as regression analysis, machine learning, and neural networks, have also been around for some time.
- Prescriptive analytics are often referred to as advanced analytics.

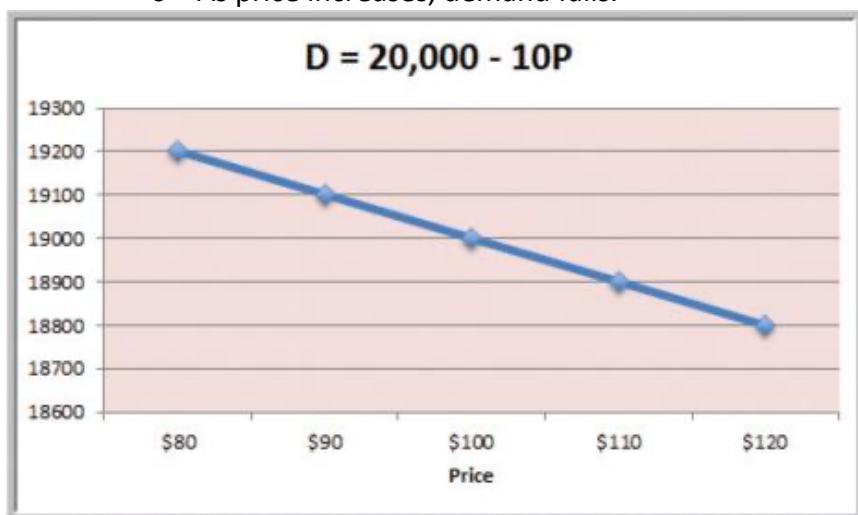


## What will occur?

- Marketing is the target for many predictive analytics applications.
- Descriptive analytics, such as data visualization, is important in helping users interpret the output from predictive and prescriptive analytics.

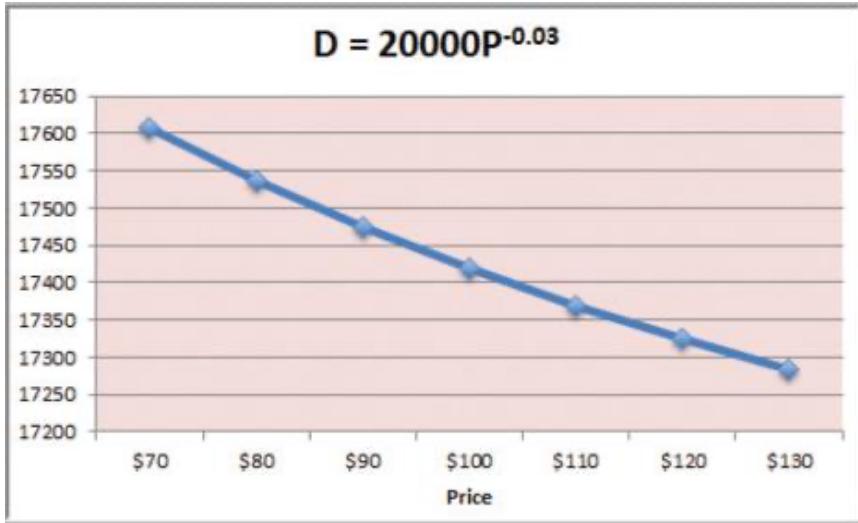
## Decision Models

- A Linear Demand Prediction Model:
  - o As price increases, demand falls.



- A Non-Linear Demand Prediction Model:

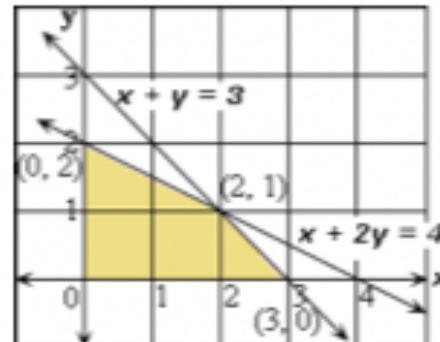
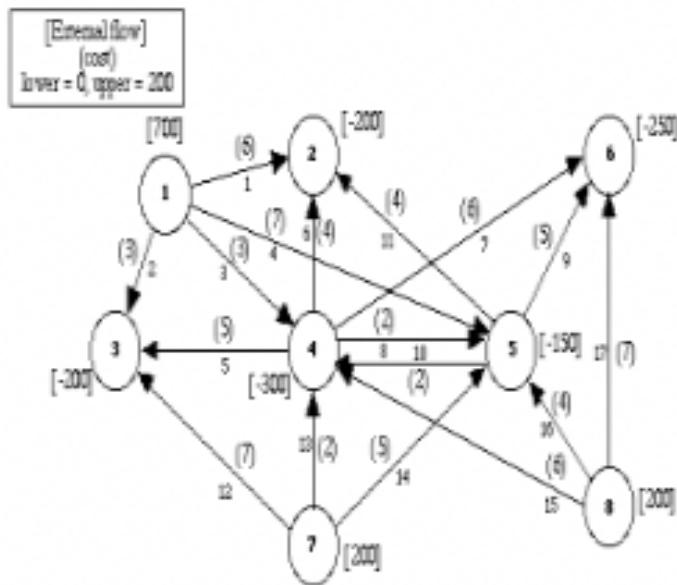
- Assumes price elasticity (constraint ratio of % change in demand to % change in price).



- **Prescriptive Decision Models** help decision makers identify the best solution.
  - Optimization – finding values for decision variables that minimize (or maximize) something such as cost (or profit).
  - Objective function – the equation that minimizes (or maximizes) the quantity of interest.
  - Constraints – limitation or restrictions.
  - Optimal solution – values of decision variables at the minimum (or maximum) point.

## Prescriptive Analytics

- Prescriptive analytics are often referred to as advanced analytics.
- Regression analysis, machine learning and neural networks.
- Often for the allocation of scarce resources.



## What should occur?

- For example, the use of mathematical programming for revenue management is common for organizations that have “perishable” good (e.g., rental cars, hotel rooms, airline seats).
- Harrah’s has been using revenue management for hotel room pricing for some time.

## Organizational Transformation

- Brought about by opportunity or necessity.
- The firm adopts a new business model enabled by analytics.
- Analytics are competitive requirements.



## 2013 Academic Research

- A 2011 TDWI report on Big Data Analytics found that 85% of respondents indicated that their firms would be using advanced analytics within 3 years.
- A 2011 IBM/MIT Sloan Management Review research study found that top performing companies in their industry are much more likely to use analytics rather than intuition across the widest range of possible decisions.

## Conditions that Lead to Analytics-based Organizations

- The nature of the industry.
- Seizing an opportunity.
- Responding to a problem.

## Complex Systems

- Tackle complex problems and provide individualized solutions.
- Products and services are organized around the needs of individual customers.
- Dollar value of interactions with each customer is high.
- There is considerable interaction with each customer.
- Examples: IBM, World Bank, Halliburton.

## Volume Operations

- Serves high-volume markets through standardized products and services.
- Each customer interaction has a low dollar value.
- Customer interactions are generally conducted through technology rather than person-to-person.
- Are likely to be analytics-based.

- Examples: Amazon, eBay, Hertz.

## The Nature of the Industry: Online Retailers

- BI Applications
  - o Analysis of clickstream data.
  - o Customer profitability analysis.
  - o Customer segmentation analysis.
  - o Product recommendations.
  - o Campaign management.
  - o Pricing.
  - o Forecasting.
  - o Dashboards.

## The Nature of the Industry

- Online retailers like Amazon.com and Overstock.com are high volume operations who rely on analytics to complete.
- When you enter their sites, a cookie is placed on your PC and all clicks are recorded.
- Based on your blocks and any search terms, recommendation engines decide what products to display.
- After you purchase an item, they have additional information that is used in marketing campaigns.
- Customer segmentation analysis is used in deciding what promotions to send you.
- How profitable you are influencing how the customer care center treats you.
- A pricing team helps set prices and decides what prices are needed to clear out merchandise.
- Forecasting models are used to decide how many items to order for inventory.
- Dashboards monitor all aspects of organizational performance.

## Analytics Help the Cincinnati Zoo Know Its Customers

- What management, organization, and technology factors were behind the Cincinnati Zoo losing opportunities to increase revenue?
- Why was replacing legacy point-of-sale systems and implementing a data warehouse essential to an information system solution?
- How did the Zoo benefit from business intelligence? How did it enhance operational performance and decision making? What role was played by predictive analytics?
- Visit the IBM Cognos Web site and describe the business intelligence tools that would be the most useful for the Zoo.

## Table of Contents

<b>What are Data Analytics? .....</b>	<b>1</b>
<b>Business Analytics applications: .....</b>	<b>1</b>
<b>Importance of Business Analytics: .....</b>	<b>1</b>
<b>Scope of Business Analytics.....</b>	<b>1</b>
<b>Scope of Business Analytics - Retail Markdown Decisions.....</b>	<b>1</b>
<b>Data for Business Analytics.....</b>	<b>2</b>
<b>Data for Business Analytics.....</b>	<b>2</b>

<i>Data for Business Analytics</i> .....	2
<i>What is Big Data?</i> .....	2
<i>Types of Data</i> .....	3
<i>Data for Business Analytics</i> .....	3
<i>Categorical (Nominal) data</i> .....	3
<i>Ordinal data</i> .....	4
<i>Interval and ratio data</i> .....	5
<i>Interval data</i> .....	5
<i>Ratio data</i> .....	5
<i>Types of analytics – Decision Models</i> .....	6
<i>Descriptive Analytics</i> .....	6
<i>Decision Models</i> .....	7
<i>Predictive Decision Models</i> .....	8
<i>Predictive Analytics</i> .....	8
<i>What will occur?</i> .....	8
<i>Decision Models</i> .....	8
<i>Prescriptive Analytics</i> .....	9
<i>What should occur?</i> .....	9
<i>Organizational Transformation</i> .....	10
<i>2013 Academic Research</i> .....	10
<i>Conditions that Lead to Analytics-based Organizations</i> .....	10
<i>Complex Systems</i> .....	10
<i>Volume Operations</i> .....	10
<i>The Nature of the Industry: Online Retailers</i> .....	11
<i>The Nature of the Industry</i> .....	11
<i>Analytics Help the Cincinnati Zoo Know Its Customers</i> .....	11

# Lecture 2: OLAP, Data Warehouse and Column Store

## Why we still study OLAP/Data Warehouse in BI?

- Understand the Big Data history.
  - o How does the requirement of (big) data analytics/business intelligence evolve over the time?
  - o What are the architecture and implementation techniques being developed? Will they still be useful in Big Data?
  - o Understand their limitation and what factors have changed from 90's to now?
- NoSQL is not only SQL.
- Hive/Impala aims to provide OLAP/BI for Big Data using Hadoop.

## Highlights

- OLAP
  - o Multi-relational Data model.
  - o Operators.
  - o SQL.
- Data warehouse (architecture, issues, optimizations).
- Join Processing.
- Column Stores (Optimized for OLAP workload).

Let's get back to the root in 70's: Relational Database

## Basic Structure

- Formally, given sets  $D_1, D_2, \dots, D_n$ . A **relation**  $r$  is a subset of  $D_1 \times D_2 \times \dots \times D_n$ .
- Thus, a **relation** is a set of n-tuples  $(a_1, a_2, \dots, a_n)$  where each  $a_i$  belong to  $D_i$ .
- Example:

```
customer_name = {Jones, Smith, Curry, Lindsay}
customer_street = {Main, North, Park}
customer_city = {Harrison, Rye, Pittsfield}
```

- Then:

```
r = {
    (Jones, Main, Harrison),
    (Smith, North, Rye),
    (Curry, North, Rye),
    (Lindsay, Park, Pittsfield)
}
```

- Is a relation over:

```
customer_name, customer_street, customer_city
```

## Relation Schema

- $A_1, A_2, \dots, A_n$  are **attributes**.
  - $R = (A_1, A_2, \dots, A_n)$  is a **relation schema**.
- Example:
- o  $\text{customer\_schema} = (\text{customer\_name}, \text{customer\_street}, \text{customer\_city})$
- $r(R)$  is a relation on the relation schema  $R$ .
- Example:
- o  $\text{customer}(\text{customer\_schema})$

## Relation Instance

- The current values (**relation instance**) of a relation are specified by a table.
- An element  $t$  of  $r$  is a **tuple**, represented by a row in a table.

The diagram shows a table representing a relation instance. The table has three columns labeled *customer\_name*, *customer\_street*, and *customer\_city*. The rows contain the following data:  
customer\_name: Jones, Smith, Curry, Lindsay  
customer\_street: Main, North, North, Park  
customer\_city: Harrison, Rye, Rye, Pittsfield

Annotations point to the table:

- An arrow from the text "attributes (or columns)" points to the column headers.
- An arrow from the text "tuples (or rows)" points to the first row of data.
- The word "customer" is centered below the table.

<i>customer_name</i>	<i>customer_street</i>	<i>customer_city</i>
Jones	Main	Harrison
Smith	North	Rye
Curry	North	Rye
Lindsay	Park	Pittsfield

## Database

- A **database** consists of multiple relations.
- Information about an enterprise is broken up into parts, with each relation storing one part of the information.

```
account : stores information about accounts
depositor : stores information about which customer owns which account
customer : stores information about customers
```

- Storing all information as a single relation such as:

```
bank(account_number, balance, customer_name, ...)
```

results in **repetition** of information (e.g., two customers own an account) and the need for null values (e.g., represent a customer without an account).

## Banking Example

```
branch (branch-name, branch-city, assets)
customer (customer-name, customer-street, customer-city)
account (account-number, branch-name, balance)
loan (loan-number, branch-name, amount)
depositor (customer-name, account-number)
borrower (customer-name, loan-number)
```

## Relational Algebra

- Primitives:
  - Projection ( $\pi$ )
  - Selection ( $\sigma$ )
  - Cartesian product ( $\times$ )
  - Set union ( $\cup$ )
  - Set difference ( $-$ )
  - Rename ( $\rho$ )
- Other operations:
  - Join ( $\bowtie$ )
  - Group by... aggregation
  - ...

## What happens next?

- SQL.
- System R (DB2), INGRES, ORACLE, SQL-Server, Teradata.
  - B+-Tree (select).
  - Transaction Management.
  - Join Algorithm.

## In early 90's: OLAP & Data Warehouse

## Database Workloads

- OLTP (online transaction processing).
  - Typical applications: e-commerce, banking, airline reservations.
  - User-facing: real-time, low latency, **highly concurrent**.
  - Tasks: relatively small set of "standard" transactional queries.
  - Data access pattern: random reads, updates, writes (involving relatively small amounts of data).
- OLAP (online analytical processing).
  - Typical applications: business intelligence, data mining.
  - Back-end processing: **batch workloads, less concurrency**.
  - Tasks: complex analytical queries, often ad-hoc.
  - Data access pattern: table scans, large amounts of data involved per query.

## OLTP

- Most database operations involve (On-Line Transaction Processing).
  - o Short, simple, frequent queries and/or modifications, each involving a small number of tuples.
  - o Examples: Answering queries from Web interface, sales at cash registers, selling airline tickets.

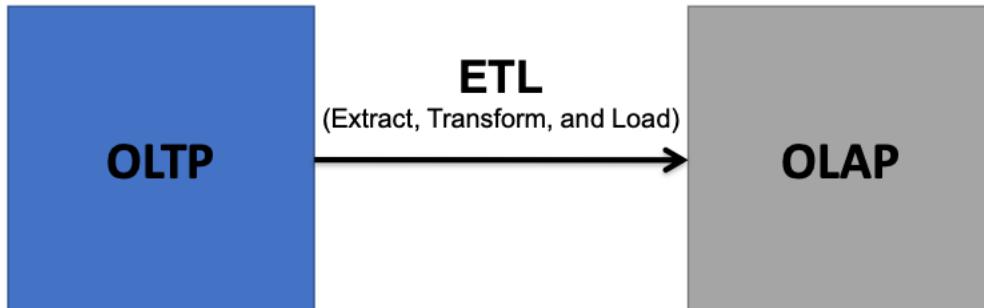
## OLAP

- Of increasing importance are On-line Application Processing (OLAP) queries.
  - o Few, but complex queries – may run for hours.
  - o Queries do not depend on having an **ABSOLUTELY** up-to-date database.
- Examples:
  - o 1. Amazon analyzes purchases by its customers to come up with an individual screen with products of likely interest to the customer.
  - o 2. Analysts at Wal-Mart look for items with increasing sales in some region.

## One Database or Two?

- Downsides of co-existing OLTP and OLAP workloads:
  - o Poor memory management.
  - o Conflicting data access patterns.
  - o Variable latency.
- Solution: separate databases.
  - o User-facing OLTP database for high-volume transaction.
  - o Data warehouse for OLAP workloads
  - o How do we connect the two?

## OLTP/OLAP Architecture



## OLTP/OLAP Integration

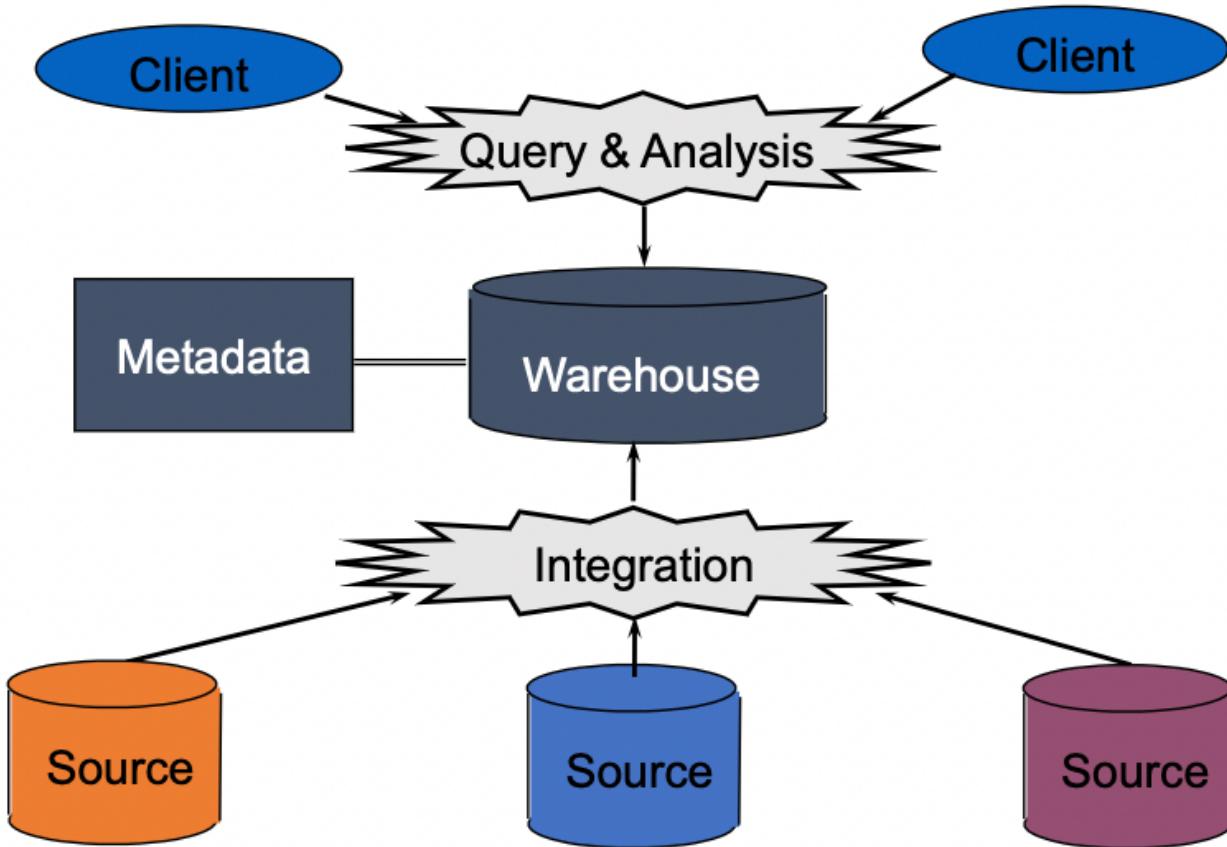
- OLTP database for user-facing transactions.
  - o Retain records of all activity.
  - o Periodic ETL (e.g., nightly).
- Extract-Transform-Load (ETL):
  - o Extract records from source.
  - o Transform: clean data, check integrity, aggregate, etc.
  - o Load into OLAP database.
- OLAP database for data warehousing:

- BI: reporting, ad-hoc queries, data mining, etc.
- Feedback to improve OLTP services.

## The Data Warehouse

- The most common form of data integration.
  - Copy sources into a single DB (warehouse) and try to keep it up to date.
  - Usual method: periodic reconstruction of the warehouse, perhaps overnight.
  - Frequently essential for analytic queries.

## Warehouse Architecture



## Star Schemas

- A **star schema** is a common organization for data at a warehouse. It consists of:
  - 1. Fact table:** a very large accumulation of facts such as sales.
    - Often “insert-only”.
  - 2. Dimension tables:** smaller, generally static information about the entities involved in the facts.

## Example: Star Schema

- Suppose we want to record in a warehouse information about every beer sale: the bar, the brand of beer, the drinker who bought the beer, the day, the time, and the price charged.
- The fact table is a relation:

**Sales(bar, beer, drinker, day, time, price)**

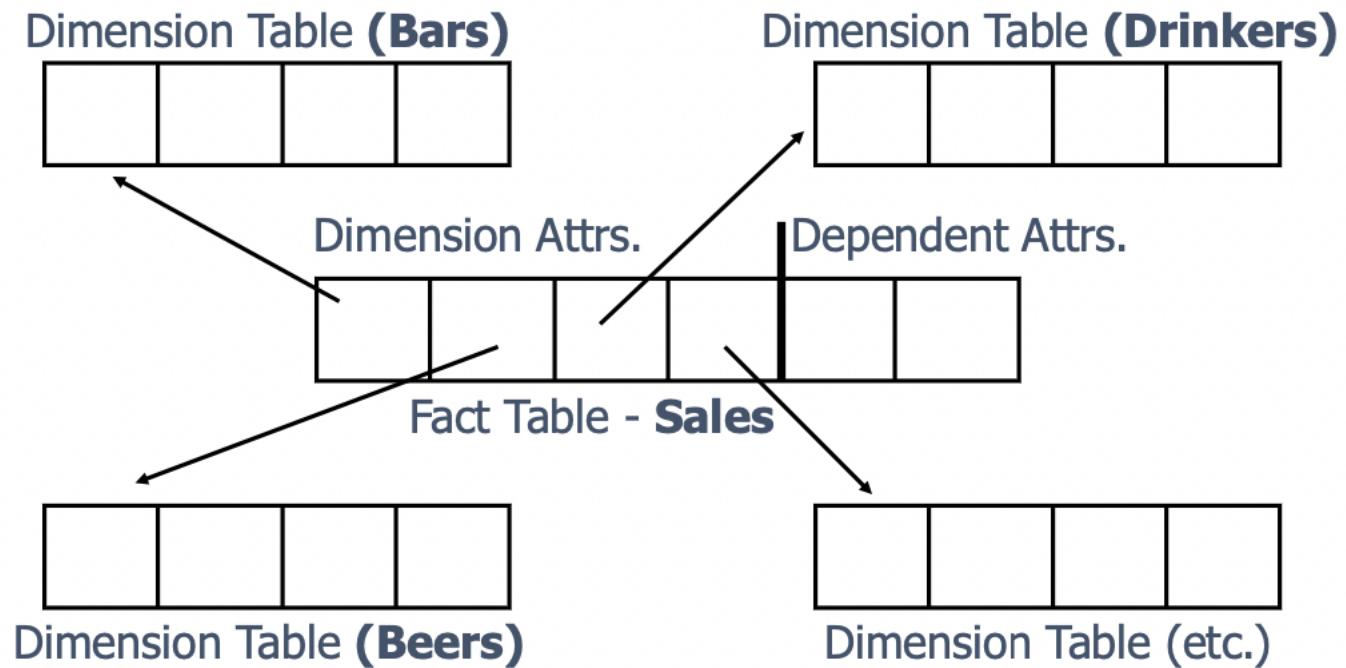
- The dimension tables include information about the bar, beer, and drinker “dimensions”:

**Bars(bar, addr, license)**

**Beers(beer, manf)**

**Drinkers(drinker, addr, phone)**

Visualization – Start Schema



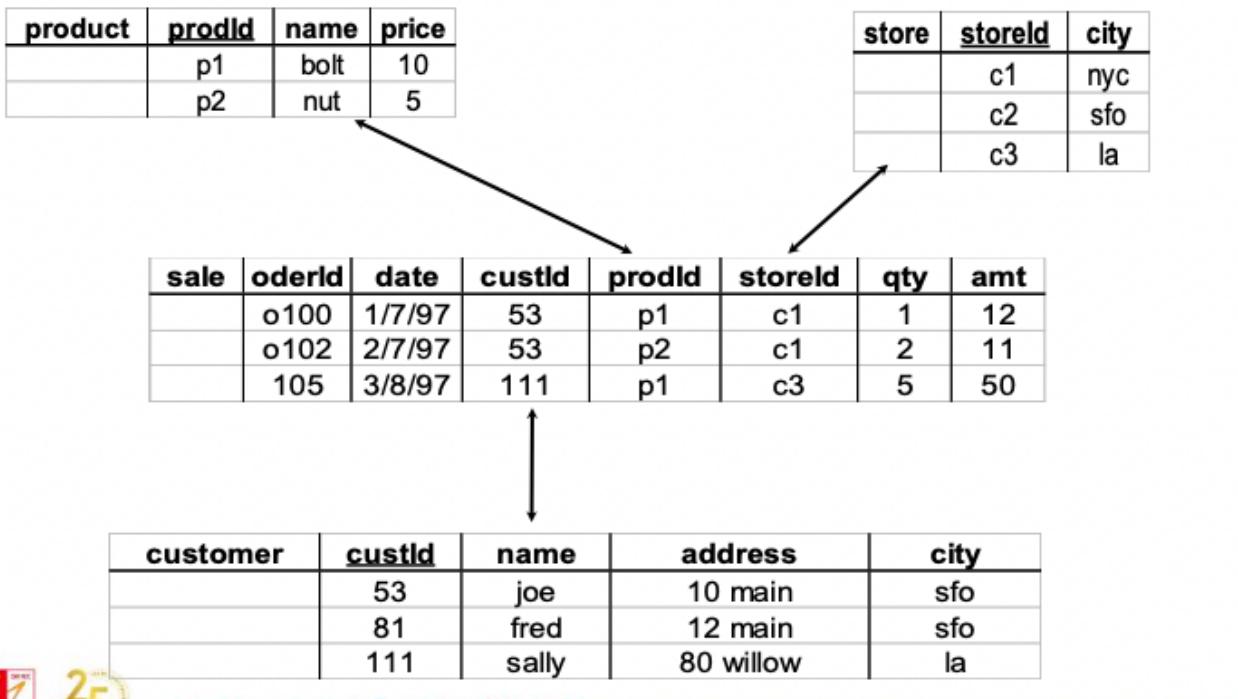
### Dimension and Dependent Attributes

- Two classes of fact-table attributes:
  - o **1. Dimension attributes:** the keys of dimension tables.
  - o **2. Dependent attributes:** a value determined by the dimension attributes of tuple.

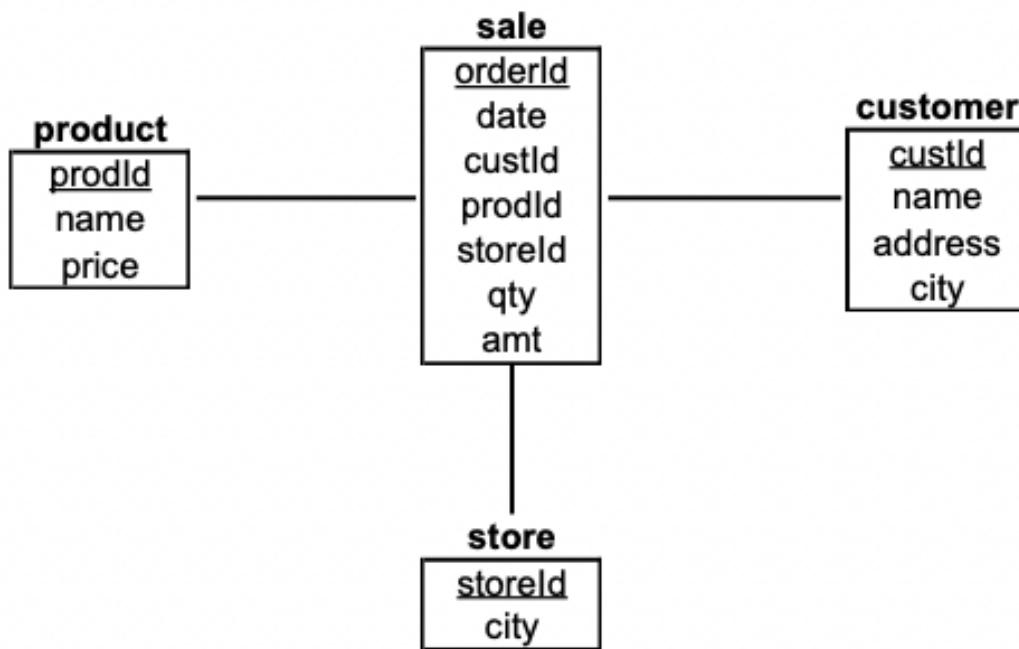
### Warehouse Models & Operators

- Data Models:
  - o Relations.
  - o Stars & snowflakes.
  - o Cubes.
- Operators:
  - o Slice & dice.
  - o Roll-up, drill-down.
  - o Pivoting.
  - o Other.

## Star



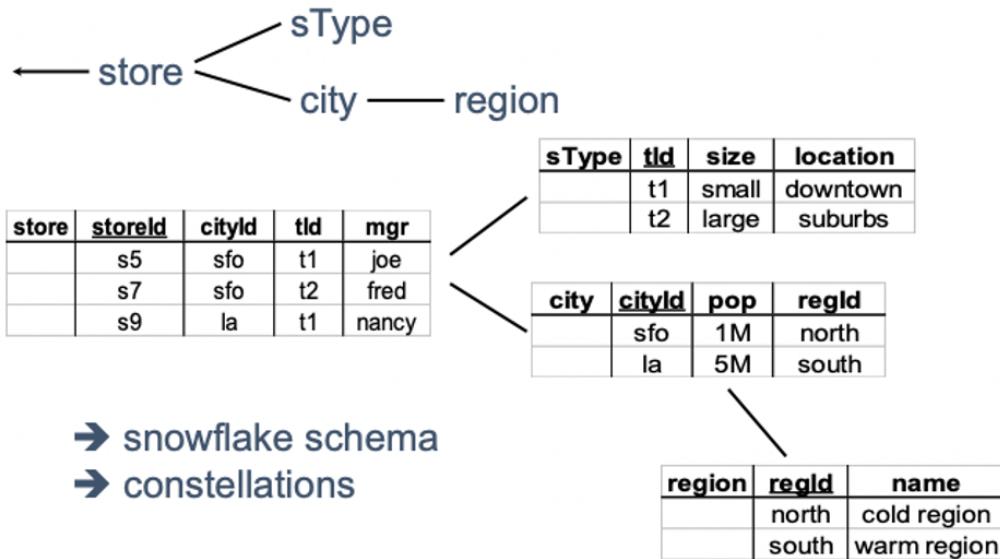
## Star Schema



## Terms

- Fact table.
- Dimension tables.
- Measures.

## Dimension Hierarchies



## Aggregates

- In SQL: **SELECT sum(amt) FROM SALE WHERE date = 1**

sale	prodId	storeId	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

→ 81

**SELECT date, sum(amt) FROM SALE GROUP BY date**

sale	prodId	storeId	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

→

ans	date	sum
	1	81
	2	48

**SQL: SELECT date, sum(amt) FROM SALE GROUP BY date, prodId**

sale	prodId	storeId	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

→

sale	prodId	date	amt
	p1	1	62
	p2	1	19
	p1	2	48

→ rollup →

← drill-down ←

## ROLAP VS. MOLAP

- ROLAP: Relational On-Line Analytical Processing
- MOLAP: Multi-Dimensional On-Line Analytical Processing

### Cube

Fact table view:

sale	prodId	storeId	amt
	p1	c1	12
	p2	c1	11
	p1	c3	50
	p2	c2	8

Multi-dimensional cube:

	c1	c2	c3
p1	12		50
p2	11	8	

dimensions = 2

### 3-D Cube

Fact table view:

sale	prodId	storeId	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

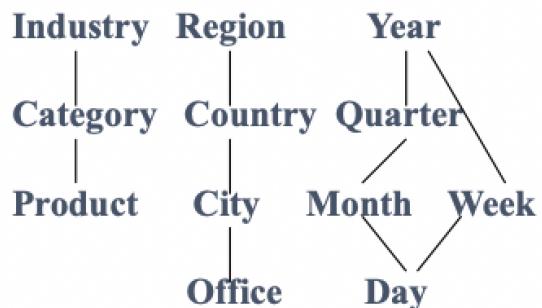
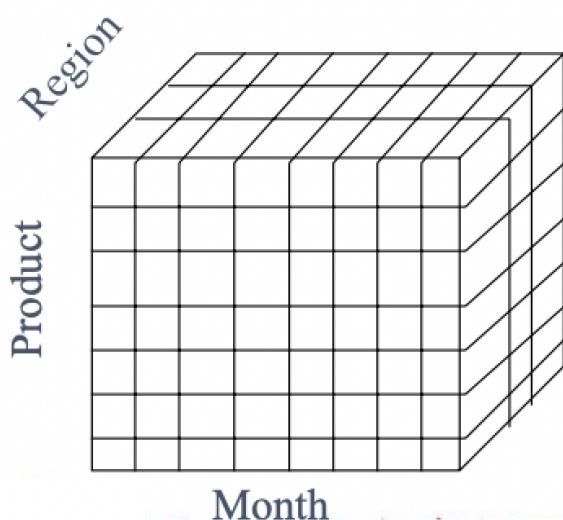
Multi-dimensional cube:

day 2		c1	c2	c3
	p1	44	4	
day 1		c1	c2	c3
	p1	12		50
	p2	11	8	

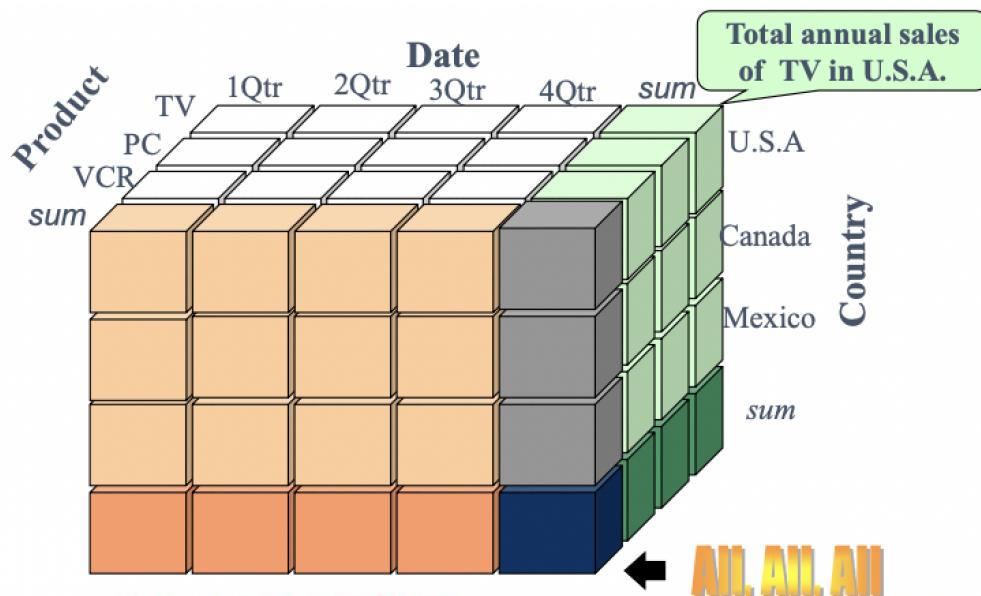
### Multidimensional Data

- Sales volume as a function of product, month, and region.

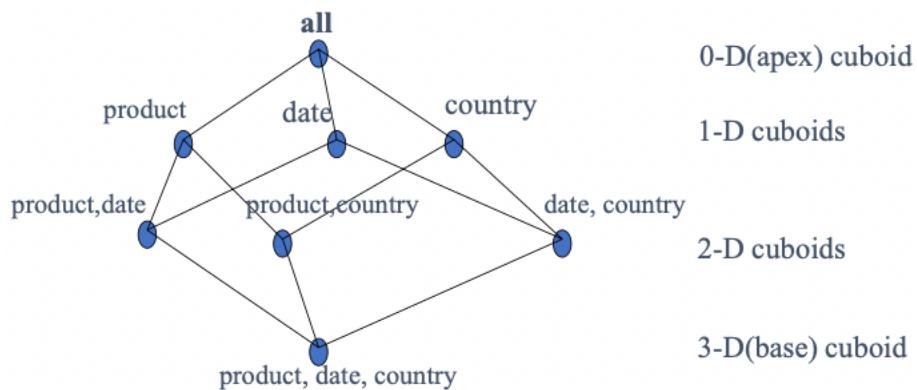
Dimensions: Product, Location, Time  
Hierarchical summarization paths



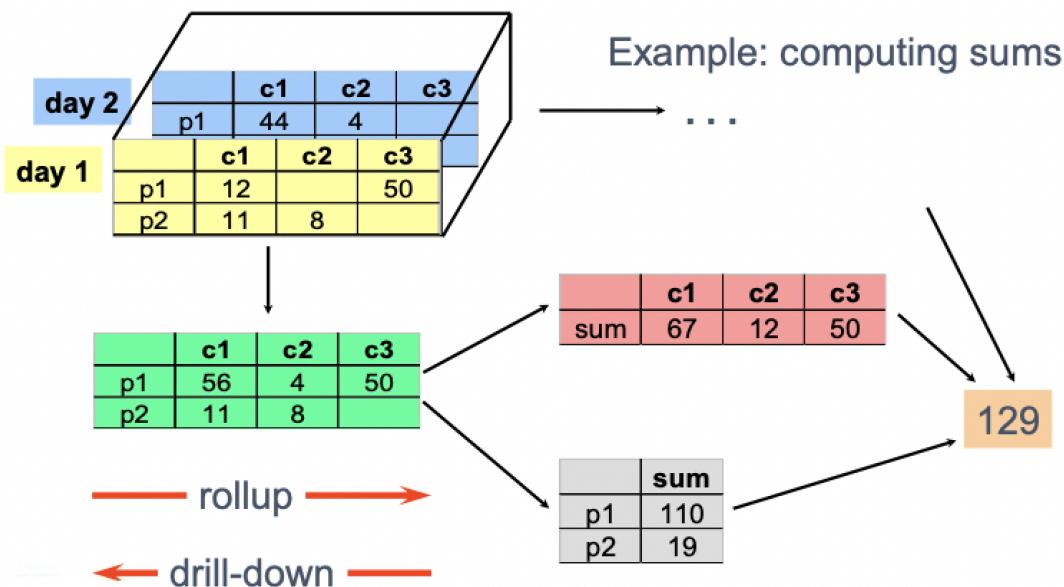
## A Sample Data Cube



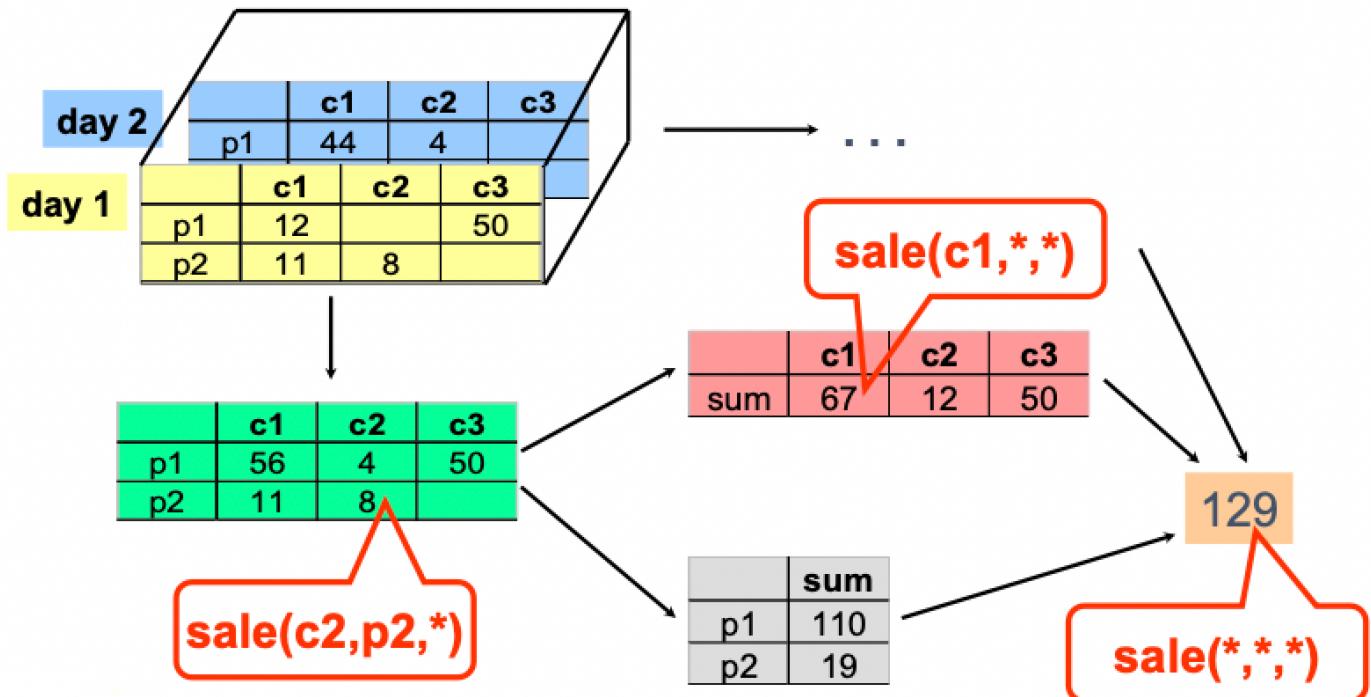
## Cuboids Corresponding to the Cube



## Cube Aggregation



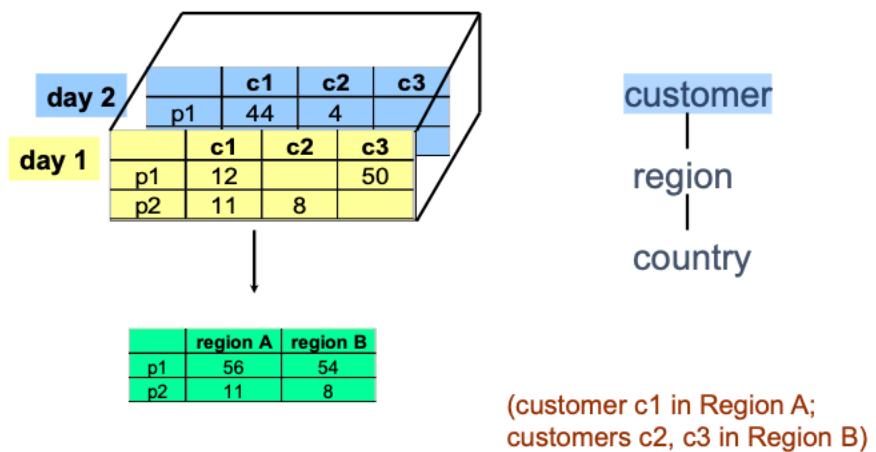
## Cube Operators



## Extended Cube

	*	c1	c2	c3	*	sum
day 2		p1	56	4	50	110
		p2	11	8		19
day 1		p1	44	4		48
		p2	11	8		19
*		*	23	8	50	81

## Aggregation Using Hierarchies

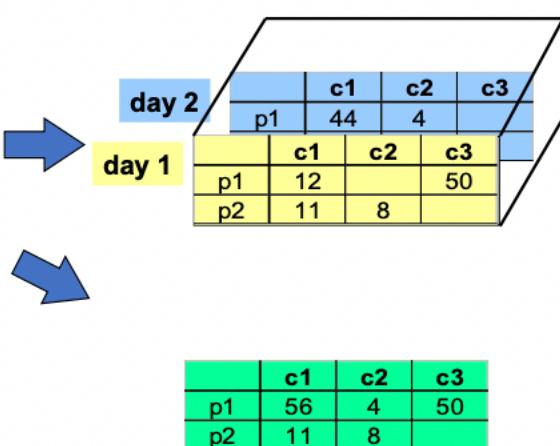


## Pivoting

Fact table view:

sale	prodId	storeId	date	amt
	p1	c1	1	12
	p2	c1	1	11
	p1	c3	1	50
	p2	c2	1	8
	p1	c1	2	44
	p1	c2	2	4

Multi-dimensional cube:



CUBE Operator (SQL-99)

Chevy Sales Cross Tab				
Chevy	1990	1991	1992	Total (ALL)
black	50	85	154	289
white	40	115	199	354
Total (ALL)	90	200	353	1286

```

SELECT model, year, color, sum(sales) as sales
FROM   sales
WHERE  model in ( 'Chevy' )
AND    year BETWEEN 1990 AND 1992
GROUP BY CUBE (model, year, color);
  
```

- Computes union of 8 different groups:
  - o {(model, year, color), (model, year), (model, color), (year, color), (model), (year), (color), ()}

## Aggregates

- Operators: sum, count, max, min, median, average.
- “Having” clause.
- Cube (& Rollup) operator.
- Using dimension hierarchy.
  - average by region (within store).
  - maximum by month (within date).

## Query & Analysis Tools

- Query Building.
- Report Writers (comparisons, growth, graphs, etc.)

- Spreadsheet Systems.
- Web Interfaces.
- Data Mining.

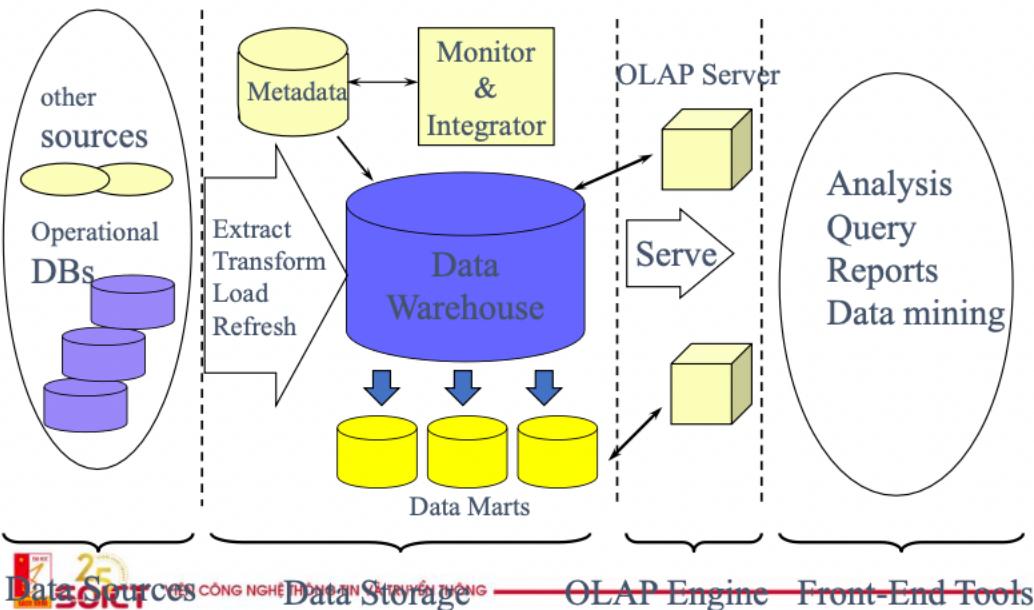
## Other Operations

- Time functions.
  - o E.g., time average.
- Computed Attributes:
  - o E.g., commission = sales \* rate.
- Text Queries:
  - o E.g., find documents with words X AND B.
  - o E.g., rank documents by frequency of words X, Y, Z.

## Data Warehouse Implementation

- Monitoring: Sending data from sources.
- Integrating: Loading, Cleansing.
- Processing: Query processing, indexing, etc.
- Managing: Metadata, Design.

## Multi-Tiered Architecture



## Monitoring

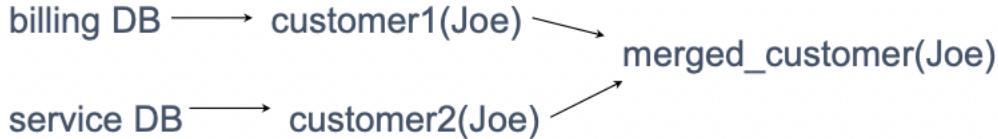
- Source Types: relational, flat file, IMS, VSAM, IDMS, WWW, newswire, etc.
- Incremental vs. Refresh.

customer	id	name	address	city
	53	joe	10 main	sfo
	81	fred	12 main	sfo
	111	sally	80 willow	la



## Data Cleansing

- Migration (e.g., yen → dollars).
- Scrubbing: use domain-specific knowledge (e.g., social security numbers).
- Fusion (e.g., mail list, customer merging).
- Auditing: discover rules & relationships (like data mining).



## Loading data

- Incremental vs. refresh.
- Off-line vs. on-line.
- Frequency of loading.
- At night, 1x a week/month, continuously.
- Parallel/Partitioned load.

## OLAP Implementation

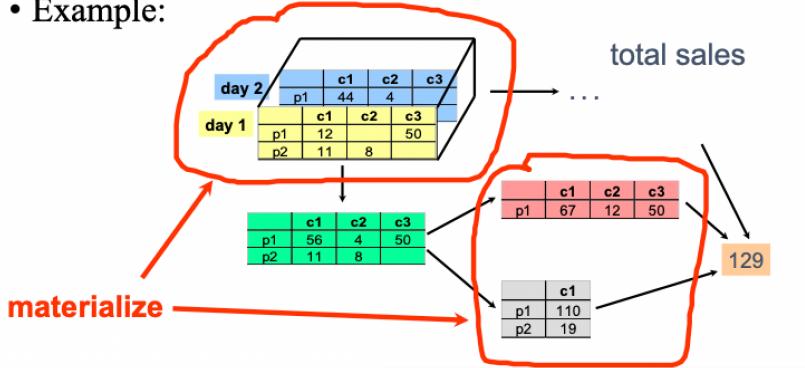
### Derived Data

- Derived Warehouse Data
  - o Indexes
  - o Aggregates
  - o materialized views
- When to update derived data?
- Incremental vs. refresh.

### What to Materialize

- Store in warehouse results useful for common queries.

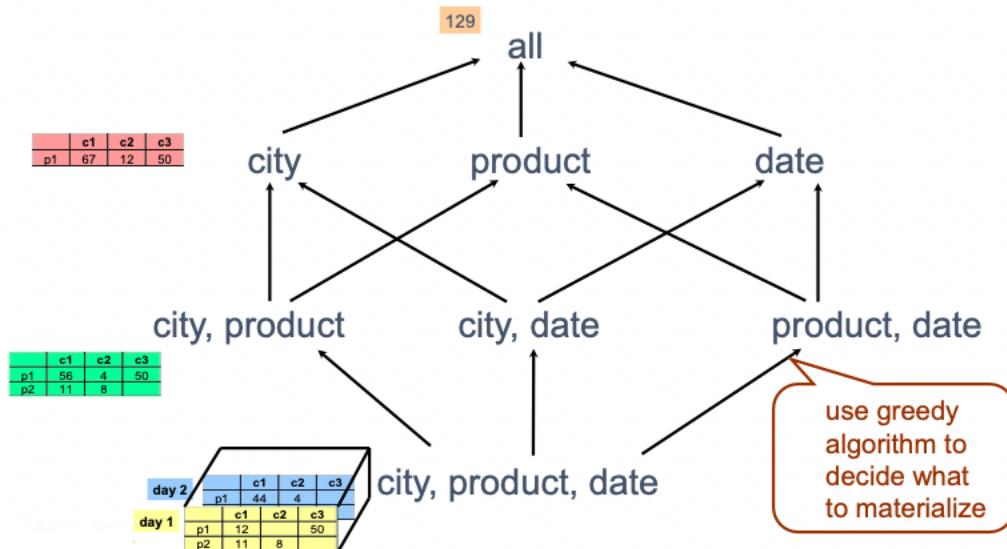
• Example:



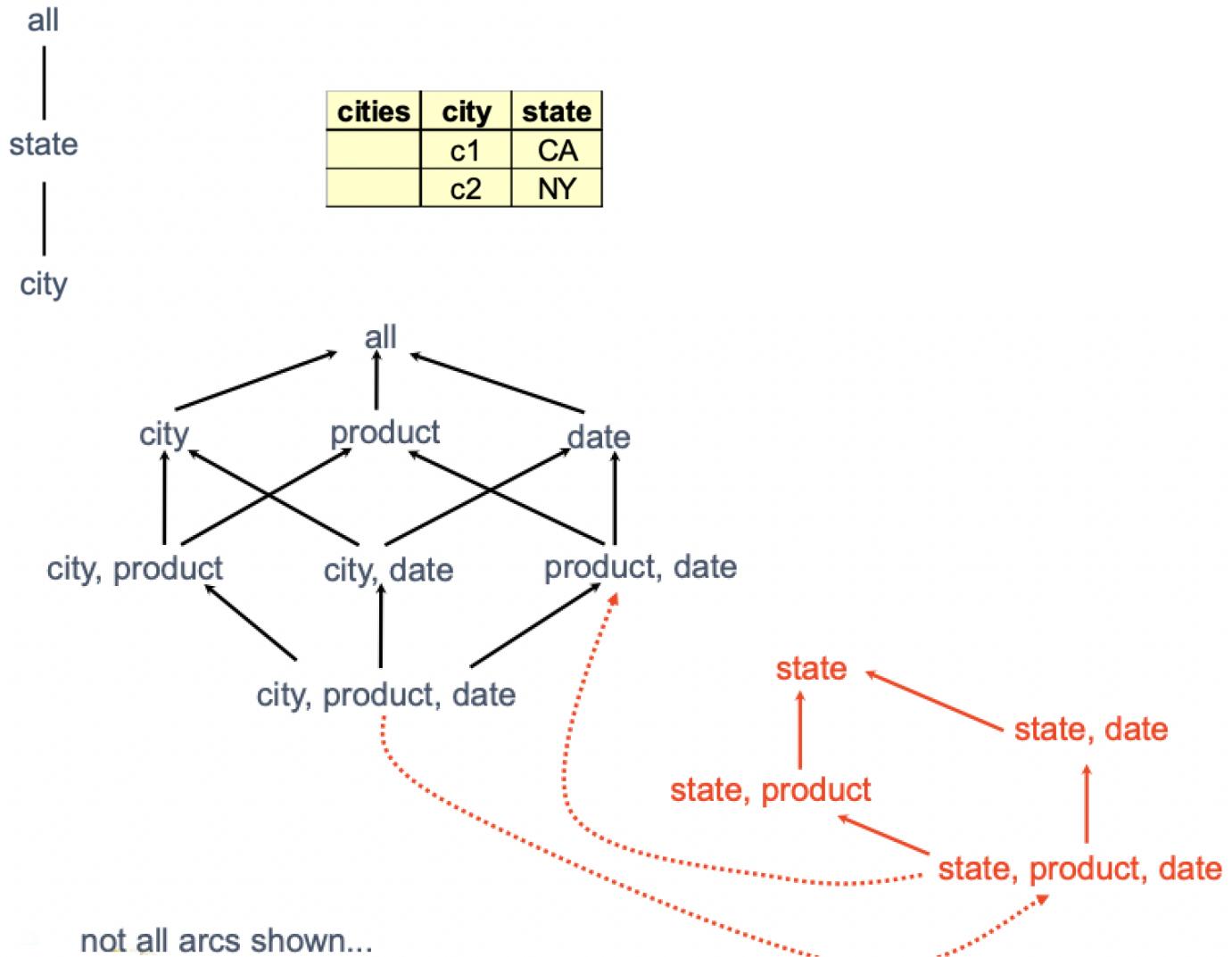
### Materialization Factors

- Type/frequency of queries.
- Query response time.
- Storage Cost.
- Update Cost.

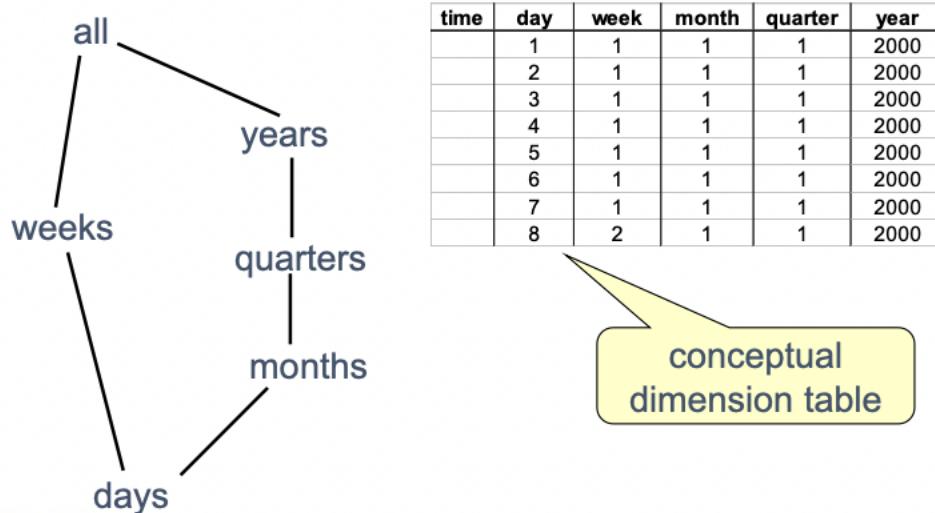
## Cube Aggregates Lattice



## Dimension Hierarchies



## Interesting Hierarchy



## Indexing OLAP Data: Bitmap Index

- Index on a particular column
- Each value in the column has a bit vector: bit-op is fast
- The length of the bit vector: # of records in the base table
- The i-th bit is set if the i-th row of the base table has the value for the indexed column
- not suitable for high cardinality domains

Base table			Index on Region			Index on Type			
Cust	Region	Type	RecID	Asia	Europe	America	RecID	Retail	Dealer
C1	Asia	Retail	1	1	0	0	1	1	0
C2	Europe	Dealer	2	0	1	0	2	0	1
C3	Asia	Dealer	3	1	0	0	3	0	1
C4	America	Retail	4	0	0	1	4	1	0
C5	Europe	Dealer	5	0	1	0	5	0	1

**SQCT** (SELECT WHEN COLUMNS COUNT IS TRUE)

## Join Processing

- How does DBMS join 2 tables?
- Sorting is one way.
- Database must choose the best way for each query.

## Schema for Examples

**Sailors** (sid: integer, sname: string, rating: integer, age: real)  
**Reserves** (sid: integer, bid: integer, day: dates, rname: string)

- Like old schema, `rname` added for variations.
- Reserves:
  - o Each tuple is 40 bytes long.
  - o 100 tuples per page.
  - o M = 1000 pages total.
- Sailors:

- Each tuple is 50 bytes long.
- 80 tuples per page.
- $N = 500$  pages total.

## Equality Joins with One Join Column

```
SELECT *
FROM   Reserves R1, Sailors S1
WHERE  R1.sid=S1.sid
```

- In algebra:  $R \bowtie S$ . Common! Must be carefully optimized.  
 $R \times S$  is large; so,  $R \times S$  followed by a selection is inefficient.
- Assume:  $M$  tuples in  $R$ ,  $p_R$  tuples per page,  $N$  tuples in  $S$ ,  $p_S$  tuples per page.
  - In our examples,  $R$  is Reserves and  $S$  is Sailors.
- We will consider more complex join conditions later.
- **Cost metric:** # of I/Os. We will ignore output costs.

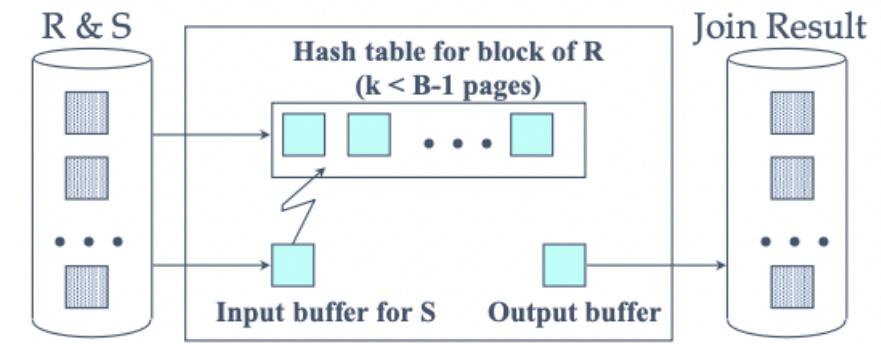
## Simple Nested Loops Join

```
foreach tuple r in R do
  foreach tuple s in S do
    if  $r_i == s_j$  then add  $\langle r, s \rangle$  to result
```

- For each tuple in the *outer* relation  $R$ , we scan the entire *inner* relation  $S$ .
  - Cost:  $M + p_R * M * N = 1000 + 100 * 1000 * 500$  I/Os.
- Page-oriented Nested Loops join: For each *page* of  $R$ , get each *page* of  $S$ , and write out matching pairs of tuples  $\langle r, s \rangle$ , where  $r$  is in  $R$ -page and  $S$  is in  $S$ -page.
  - Cost:  $M + M * N = 1000 + 1000 * 500$
  - If smaller relation ( $S$ ) is outer, cost =  $500 + 500 * 1000$

## Block Nested Loops Join

- Use one page as an input buffer for scanning the inner  $S$ , one page as the output buffer, and use all remaining pages to hold ‘‘block’’ of outer  $R$ .
- For each matching tuple  $r$  in  $R$ -Block,  $s$  in  $S$ -page, add  $\langle r, s \rangle$  to result. Then read next  $R$ -block, scan  $S$ , etc.



## Examples of Block Nested Loops

- Cost: Scan of outer + #outer blocks \* scan of inner
  - #outer blocks =  $\lceil \# \text{ of pages of outer} / \text{blocksize} \rceil$
- With Reserves (R) as outer, and 100 pages of R:
  - Cost of scanning R is 1000 I/Os; a total of 10 blocks.
  - Per block of R, we scan Sailors (S); 10\*500 I/Os.
  - If space for just 90 pages of R, we would scan S 12 times.
- With 100-page block of Sailors as outer:
  - Cost of scanning S is 500 I/Os; a total of 5 blocks.
  - Per block of S, we scan Reserves; 5\*1000 I/Os.
- With sequential reads considered, analysis changes: may be best to divide buffers evenly between R and S.

## Index Nested Loops Join

```
foreach tuple r in R do
    foreach tuple s in S where  $r_i == s_j$  do
        add  $\langle r, s \rangle$  to result
```

- If there is an index on the join column of one relation (say S), can make it the inner and exploit the index.
  - Cost:  $M + (M * p_R) * \text{cost of finding matching S tuples}$
- For each R tuple, cost of probing S index is about 1.2 for hash index, 2-4 for B+ tree. Cost of then finding S tuples (assuming Alt. (2) or (3) for data entries) depends on clustering.
  - Clustered index: 1 I/O (typical), unclustered: upto 1 I/O per matching S tuple.

## Examples of Index Nested Loops

- Hash-index (Alt. 2) on *sid* of Sailors (as inner):
  - Scan Reserves: 1000 page I/Os, 100\*1000 tuples.
  - For each Reserves tuple: 1.2 I/Os to get data entry in index, plus 1 I/O to get (the exactly one) matching Sailors tuple.  
Total: 220,000 I/Os.
- Hash-index (Alt. 2) on *sid* of Reserves (as inner):
  - Scan Sailors: 500 page I/Os, 80\*500 tuples.
  - For each Sailors tuple: 1.2 I/Os to find index page with data entries, plus cost of retrieving matching Reserves tuples.  
Assuming uniform distribution, 2.5 reservations per sailor (100,000 / 40,000). Cost of retrieving them is 1 or 2.5 I/Os depending on whether the index is clustered.

## Sort-Merge Join

$$(R \bowtie S)_{i=j}$$

- Sort R and S on the join column, then scan them to do a ``merge`` (on join column), and input result tuples.
  - o Advance scan of R until current R-tuple  $\geq$  current S tuple, then advance scan of S until current S-tuple  $\geq$  current R tuple; do this until current R tuple = current S tuple.
  - o At this point, all R tuples with same value in  $R_i$  (current R group) and all S tuples with same value in  $S_j$  (current S group) match; output  $\langle r, s \rangle$  for all pairs of such tuples.
  - o • Then resume scanning R and S.
- R is scanned once; each S group is scanned once per matching R tuple. (Multiple scans of an S group are likely to find needed pages in buffer).

### Example of Sort-Merge Join

		<u>sid</u>	<u>bid</u>	<u>day</u>	<u>rname</u>
<u>sid</u>	<u>sname</u>	<u>rating</u>	<u>age</u>		
22	dustin	7	45.0	28	103 12/4/96 guppy
28	yuppy	9	35.0	28	103 11/3/96 yuppy
31	lubber	8	55.5	31	101 10/10/96 dustin
44	guppy	5	35.0	31	102 10/12/96 lubber
58	rusty	10	35.0	31	101 10/11/96 lubber
				58	103 11/12/96 dustin

- **Cost:**  $M \log M + N \log N + (M+N)$ 
  - The cost of scanning,  $M+N$ , could be  $M*N$  (very unlikely!)
- With 35, 100 or 300 buffer pages, both Reserves and Sailors can be sorted in 2 passes; total join cost: 7500.

### Refinement of Sort-Merge Join

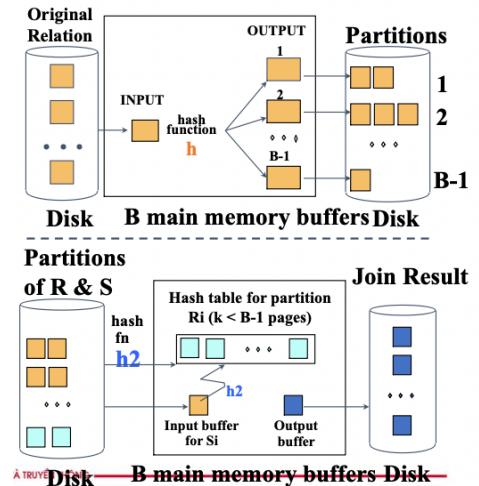
- We can combine the merging phases in the *sorting* of R and S with the merging required for the join.
  - With  $B > \sqrt{L}$ , where  $L$  is the size of the larger relation, using the sorting refinement that produces runs of length  $2B$  in Pass 0, #runs of each relation is  $< B/2$ .
  - Allocate 1 page per run of each relation, and 'merge' while checking the join condition.
  - **Cost:** read+write each relation in Pass 0 + read each relation in (only) merging pass (+ writing of result tuples).
  - In example, cost goes down from 7500 to 4500 I/Os.
- In practice, cost of sort-merge join, like the cost of external sorting, is *linear*.

## Hash-Join

- Partition both relations using hash fn  $h$ : R tuples in partition i will only match S tuples in partition i.
- Read in a partition of R, hash it using  $h_2 (< h!)$ . Scan matching partition of S, search for matches.

## Observations in Hash-Join

- #partitions  $k < B-1$  (why?), and  $B-2 > \text{size of largest partition}$  to be held in memory. Assuming uniformly sized partitions, and maximizing k, we get:
  - $k = B-1$ , and  $M/(B-1) < B-2$ , i.e., B must be  $> \sqrt{M}$
- If we build an in-memory hash table to speed up the matching of tuples, a little more memory is needed.
- If the hash function does not partition uniformly, one or more R partitions may not fit in memory. Can apply hash-join technique recursively to do the join of this R-partition with corresponding S-partition.

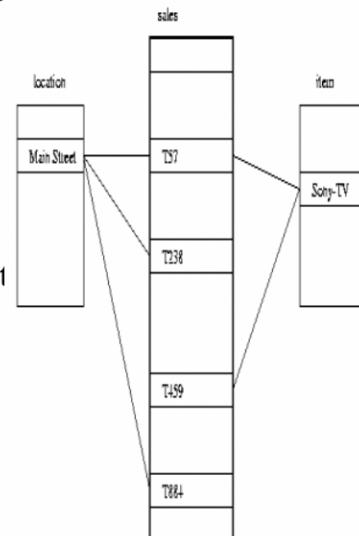


## Cost of Hash-Join

- In partitioning phase, read+write both relns;  $2(M+N)$ . In matching phase, read both relns;  $M+N$  I/Os.
- In our running example, this is a total of 4500 I/Os.
- Sort-Merge Join vs. Hash Join:
  - Given a minimum amount of memory (*what is this, for each?*) both have a cost of  $3(M+N)$  I/Os. Hash Join superior on this count if relation sizes differ greatly. Also, Hash Join shown to be highly parallelizable.
  - Sort-Merge less sensitive to data skew; result is sorted.

## Join Indices

- Traditional indices map the values to a list of record ids
  - It materializes relational join in JI file and speeds up relational join — a rather costly operation
- In data warehouses, join index relates the values of the dimensions of a star schema to rows in the fact table.
  - E.g. fact table: *Sales* and two dimensions *city* and *product*
    - A join index on *city* maintains for each distinct city a list of R-IDs of the tuples recording the Sales in the city
  - Join indices can span multiple dimensions



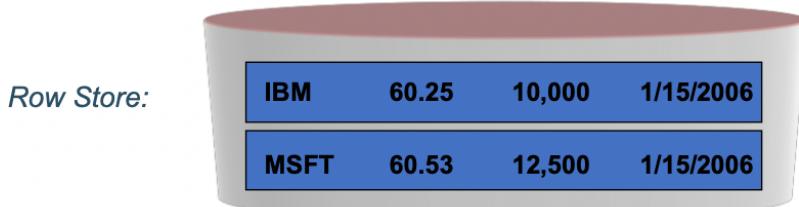
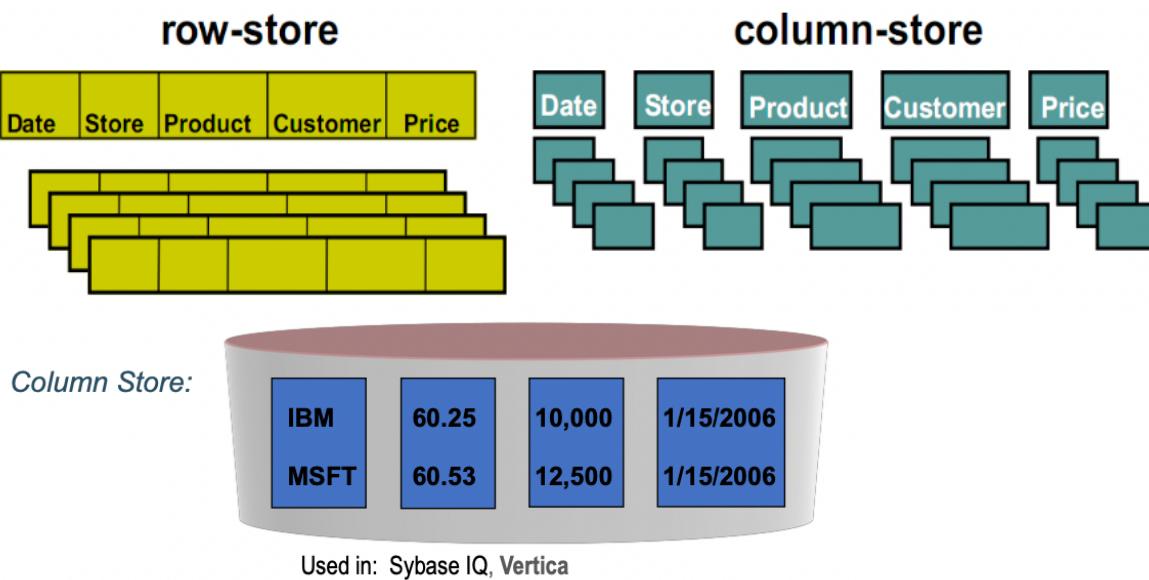
## General Join Conditions

- Equalities over several attributes (e.g.,  $R.sid=S.sid$  AND  $R.rname=S.sname$ ):
  - For Index NL, build index on  $\langle sid, sname \rangle$  (if S is inner); or use existing indexes on  $sid$  or  $sname$ .
  - For Sort-Merge and Hash Join, sort/partition on combination of the two join columns.
- Inequality conditions (e.g.,  $R.rname < S.sname$ ):
  - For Index NL, need (clustered!) B+ tree index.
    - Range probes on inner; # matches likely to be much higher than for equality joins.
  - Hash Join, Sort Merge Join not applicable.
  - Block NL quite likely to be the best join method here.

An invention in 2000s: Column Stores for OLAP

## Row Store and Column Store

- In row store data are stored in the disk tuple by tuple.
- Where in column store data are stored in the disk column by column



For example the query

```
SELECT account.account_number,
       sum(usage.toll_airtime),
       sum(usage.toll_price)
  FROM usage, toll, source, account
 WHERE usage.toll_id = toll.toll_id
   AND usage.source_id = source.source_id
   AND usage.account_id = account.account_id
   AND toll.type_ind in ('AE', 'AA')
   AND usage.toll_price > 0
   AND source.type != 'CIBER'
   AND toll.rating_method = 'IS'
   AND usage.invoice_date = 20051013
  GROUP BY account.account_number
```

Row-store: one row = 212 columns!

Column-store: 7 attributes

Row Store	Column Store
(+) Easy to add/modify a record	(+) Only need to read in relevant data
(-) Might read in unnecessary data	(-) Tuple writes require multiple accesses

- So, column storages are suitable for read-mostly, read intensive, large data repositories.

## Column Stores: High Level

- Read only what you need:
  - o “Fat” fact tables are typical.
  - o Analytics read only a few columns.
- Better compression.
- Execute and compressed data.
- Materialized views help row stores and column stores about equally.

## Data Model (Vertica/C-Store)

- Same as relational data model:
  - o Tables, rows, columns.
  - o Primary keys and foreign keys.
  - o **Projections:**
    - From single table.
    - Multiple joined tables.
- Example:

## Possible C-store model

### Normal relational model

**EMP**(name, age, dept, salary)

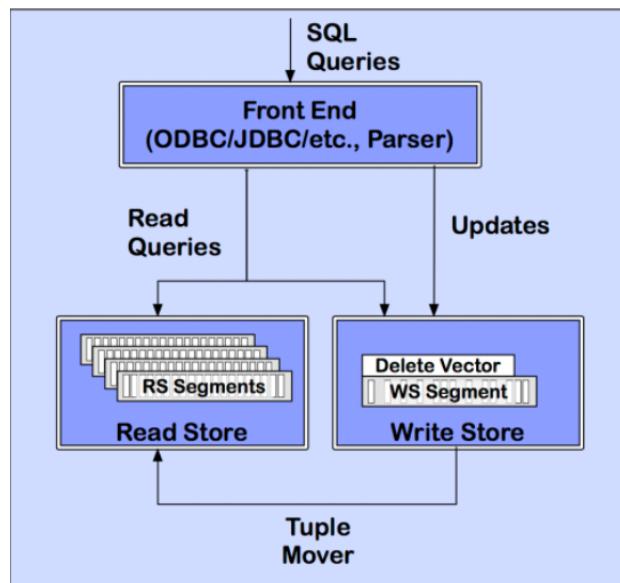
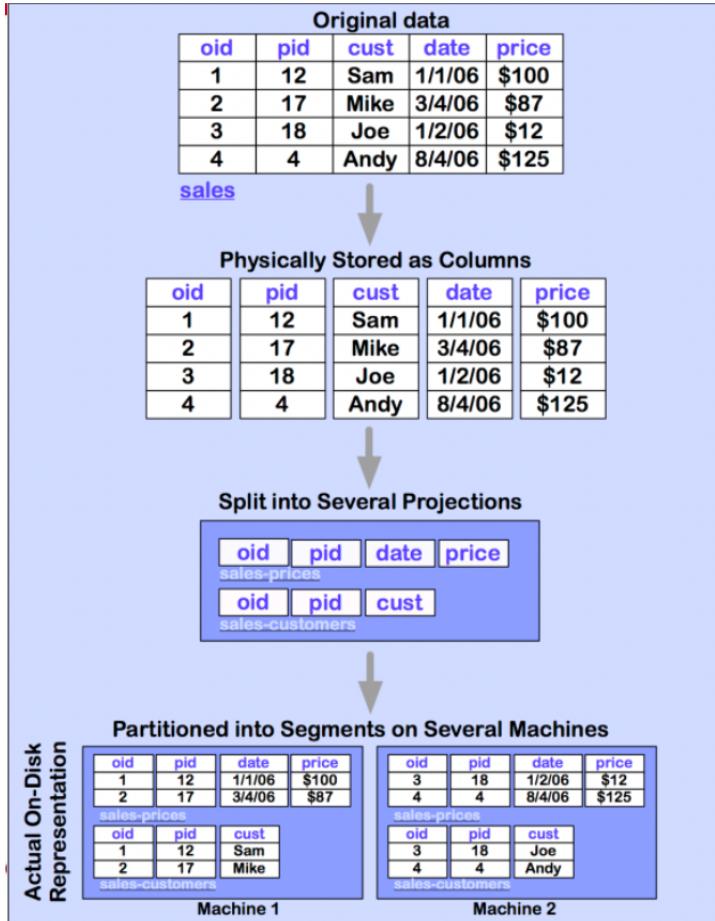
**DEPT**(dname, floor)

**EMP1** (name, age)

**EMP2** (dept, age, DEPT.floor)

**EMP3** (name, salary)

**DEPT1**(dname, floor)



### Read store: Column Encoding/Compression

- Use compression schemes and indices:
  - o Null Suppression.
  - o Dictionary encoding.
  - o Run Length encoding.
  - o Bit-Vector encoding.
  - o Self-order (key), few distinct values
    - (value, position, # items)
    - Indexed by clustered B-tree
  - o Foreign-order (non-key), few distinct values
    - (value, bitmap index)
    - B-tree index: position → values
  - o Self-order, many distinct values
    - Delta from the previous value
    - B-tree index
  - o Foreign-order, many distinct values
    - Unencoded

## Compression

- Trades I/O for CPU
  - o Increased column-store opportunities.
  - o Higher data value locality in column stores.
  - o Techniques such as run length encoding far more useful.

## Write Store

- Same structure, but explicitly use (segment, key) to identify records.
  - o Easier to maintain the mapping.
  - o Only concerns the inserted records.
- Tuple mover:
  - o Copies batch of records to RS.
- Delete record:
  - o Mark it on RS.
  - o Purged by tuple mover.

## How to solve read/write conflict

- Situation: one transaction updates the record X, while another transaction reads X.
- Use snapshot isolation.

## Query Execution – Operators

- Select: Same as relational algebra but produces a bit string.
- Project: Same as relational algebra.
- Join: Joins projections according to predicates.
- Aggregation: SQL like aggregates.
- Sort: Sort all columns of a projection.
- Decompress: Converts compressed column to uncompressed representation.
- Mask(Bitstring B, Projection Cs) => emit only those values whose corresponding bits are 1.
- Concat: Combines one or more projections sorted in the same order into a single projection.
- Permute: Permutes a projection according to the ordering defined by a join index.
- Bitstring operators: Band – Bitwise AND, Bor – Bitwise OR, Bnot – complement.

## Benefits in Query Processing

- Selection – has more indices to use.
- Projection – some “projections” already defined.
- Join – some projections are materialized joins.
- Aggregations – works on required columns only.

## Evaluation

- Use TPC-H – decision support queries.
- Storage

C-Store	Row Store	Column Store
1.987 GB	4.480 GB	2.650 GB

## Query Performance

Query	C-Store	Row Store	Column Store
Q1	0.03	6.80	2.24
Q2	0.36	1.09	0.83
Q3	4.90	93.26	29.54
Q4	2.09	722.90	22.23
Q5	0.31	116.56	0.93
Q6	8.50	652.90	32.83
Q7	2.54	265.80	33.24

- Row store uses materialized views

Query	C-Store	Row Store	Column Store
Q1	0.03	0.22	2.34
Q2	0.36	0.81	0.83
Q3	4.90	49.38	29.10
Q4	2.09	21.76	22.23
Q5	0.31	0.70	0.63
Q6	8.50	47.38	25.46
Q7	2.54	18.47	6.28

Summary: The performance gain

- Column representation: avoid reads of unused attributes.
- Storing overlapping projections – multiple orderings of a column, more choices for query optimization.
- Compression of data – more orderings of a column in the same amount of space.
- Query operators operate on compressed representation.

## Google's Dremel/Big Query: Interactive Analysis of Web-Scale Datasets

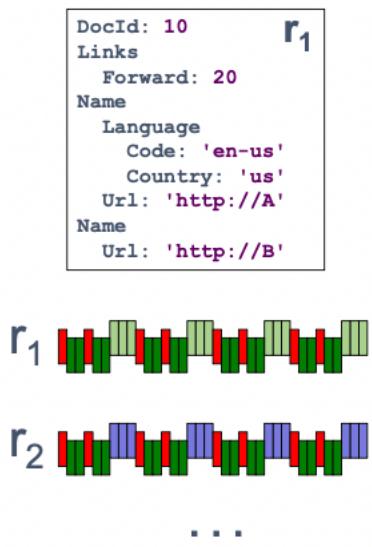
### Big Query System

- Trillion-record, multi-terabyte datasets at interactive speed
  - Scales to thousands of nodes
  - Fault and straggler tolerant execution
- Nested data model
  - Complex datasets; normalization is prohibitive
  - Columnar storage and processing
- Tree architecture (as in web search)
- Interoperates with Google's data mgmt tools
  - In situ data access (e.g., GFS, Bigtable)
  - MapReduce pipelines

## Widely Used inside Google

- Analysis of crawled web documents
- Tracking install data for applications on Android Market
- Crash reporting for Google products
- OCR results from Google Books
- Spam analysis
- Debugging of map tiles on Google Maps
- Tablet migrations in managed Bigtable instances
- Results of tests run on Google's distributed build system
- Disk I/O statistics for hundreds of thousands of disks
- Resource monitoring for jobs run in Google's data centers
- Symbols and dependencies in Google's codebase

## Records      vs.      columns



Challenge: Preserve structure, reconstruct from a subset of fields.

## Nested Data Model

<http://code.google.com/apis/protocolbuffers>

```
multiplicity:
message Document {
    required int64 DocId;           [1,1]
    optional group Links {
        repeated int64 Backward;    [0,*]
        repeated int64 Forward;
    }
    repeated group Name {
        repeated group Language {
            required string Code;
            optional string Country; [0,1]
        }
        optional string Url;
    }
}
```

DocId: 10	r <sub>1</sub>
Links	
Forward: 20	
Forward: 40	
Forward: 60	
Name	
Language	
Code: 'en-us'	
Country: 'us'	
Language	
Code: 'en'	
Url: 'http://A'	
Name	
Url: 'http://B'	
Name	
Language	
Code: 'en-gb'	
Country: 'gb'	

DocId: 20	r <sub>2</sub>
Links	
Backward: 10	
Backward: 30	
Forward: 80	
Name	
Url: 'http://C'	



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

## Column-striped representation

DocId			Name.Url			Links.Forward			Links.Backward		
value	r	d	value	r	d				value	r	d
10	0	0	http://A	0	2				N20L	0	2
20	0	0	http://B	1	2				40	0	2
			NULL	1	1				80	1	2
			http://C	0	2				80	0	2

Name.Language.Code			Name.Language.Country		
value	r	d	value	r	d
en-us	0	2	us	0	3
en	2	2	NULL	2	2
NULL	1	1	NULL	1	1
en-gb	1	2	gb	1	3
NULL	0	1	NULL	0	1



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

## Repetition and definition levels

r=1    r=2    (non-repeating)

### Name.Language.Code

value	r	d
en-us	0	2
en	2	2
NULL	1	1
en-gb	1	2
NULL	0	1

record (r=0) has repeated

Language (r=2) has repeated

```

DocId: 10          r1
Links
  Forward: 20
  Forward: 40
  Forward: 60
Name
  Language
    Code: 'en-us'
    Country: 'us'
  Language
    Code: 'en'
    Url: 'http://A'
Name
  Url: 'http://B'
Name
  Language
    Code: 'en-gb'
    Country: 'gb'
```

r: At what repeated field in the field's path  
the value has repeated

d: How many fields in paths that could be  
undefined (opt. or rep.) are actually present

```

DocId: 20          r2
Links
  Backward: 10
  Backward: 30
  Forward: 80
Name
  Url: 'http://C'
```

## Query processing

- Optimized for select-project-aggregate
  - Very common class of interactive queries
  - Single scan
  - Within-record and cross-record aggregation
- Approximations: count(distinct), top-k.
- Joins, temp tables, UDFs/TVFs, etc.

## SQL dialect for nested data

```

SELECT DocId AS Id,
       COUNT(Name.Language.Code) WITHIN Name AS Cnt,
       Name.Url + ',' + Name.Language.Code AS Str
  FROM t
 WHERE REGEXP(Name.Url, '^http') AND DocId < 20;
```

### Output table

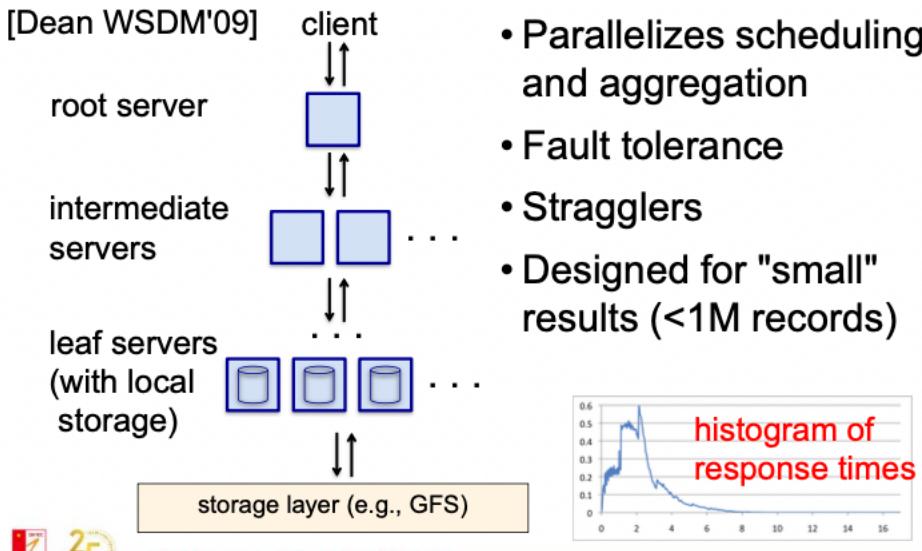
Id: 10	t <sub>1</sub>
Name	
Cnt: 2	
Language	
Str: 'http://A,en-us'	
Str: 'http://A,en'	
Name	
Cnt: 0	

### Output schema

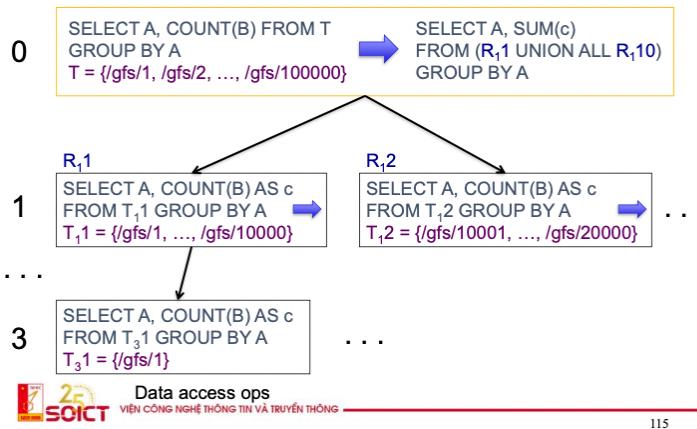
```

message QueryResult {
  required int64 Id;
  repeated group Name {
    optional uint64 Cnt;
    repeated group Language {
      optional string Str;
    }
  }
}
```

## Serving tree



## Example: count()

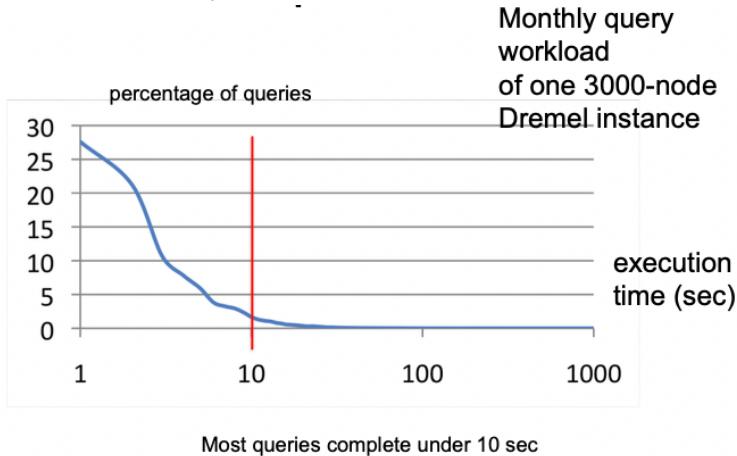


## Experiments

- 1 PB of real data (uncompressed, non-replicated).
- 100k-800k tablets per table.
- Experiments run during business hours:

Table name	Number of records	Size (unrepl., compressed)	Number of fields	Data center	Repl. factor
T1	85 billion	87 TB	270	A	3 ×
T2	24 billion	13 TB	530	A	3 ×
T3	4 billion	70 TB	1200	A	3 ×
T4	1+ trillion	105 TB	50	B	3 ×
T5	1+ trillion	20 TB	30	B	2 ×

## Interactive Speed



## Big Query: Powered by Dremel

<http://code.google.com/apis/bigquery/>



## List of Column Data Bases

- Vertica/C-Store
- SybaseIQ
- MonetDB
- LucidDB
- HANA
- Google's Dremel
- Parcell -> Redshit (Another Cloud-DB Service)

## Take-home messages

- OLAP
  - Multi-relational Data model
  - Operators
  - SQL
- Data warehouse (architecture, issues, optimizations)
- Join Processing
- Column Stores (Optimized for OLAP workload)

## Table of Contents

<b>Why we still study OLAP/Data Warehouse in BI?.....</b>	<b>1</b>
<b>Highlights.....</b>	<b>1</b>
<b>Let's get back to the root in 70's: Relational Database .....</b>	<b>1</b>
<b>Basic Structure .....</b>	<b>1</b>
<b>Relation Schema .....</b>	<b>2</b>
<b>Relation Instance .....</b>	<b>2</b>
<b>Database .....</b>	<b>2</b>
<b>Banking Example.....</b>	<b>3</b>
<b>Relational Algebra .....</b>	<b>3</b>
<b>What happens next?.....</b>	<b>3</b>
<b>In early 90's: OLAP &amp; Data Warehouse.....</b>	<b>3</b>
<b>Database Workloads.....</b>	<b>3</b>
<b>OLTP .....</b>	<b>4</b>
<b>OLAP .....</b>	<b>4</b>
<b>One Database or Two? .....</b>	<b>4</b>
<b>OLTP/OLAP Architecture .....</b>	<b>4</b>
<b>OLTP/OLAP Integration .....</b>	<b>4</b>
<b>The Data Warehouse.....</b>	<b>5</b>
<b>Warehouse Architecture.....</b>	<b>5</b>
<b>Star Schemas.....</b>	<b>5</b>
<b>Example: Star Schema .....</b>	<b>5</b>
<b>Visualization – Star Schema .....</b>	<b>6</b>
<b>Dimension and Dependent Attributes.....</b>	<b>6</b>
<b>Warehouse Models &amp; Operators .....</b>	<b>6</b>
<b>Star .....</b>	<b>7</b>
<b>Star Schema .....</b>	<b>7</b>
<b>Terms.....</b>	<b>7</b>
<b>Dimension Hierarchies.....</b>	<b>8</b>
<b>Aggregates .....</b>	<b>8</b>
<b>ROLAP VS. MOLAP.....</b>	<b>9</b>
<b>Cube.....</b>	<b>9</b>
<b>3-D Cube .....</b>	<b>9</b>
<b>Multidimensional Data.....</b>	<b>9</b>
<b>A Sample Data Cube.....</b>	<b>10</b>
<b>Cuboids Corresponding to the Cube .....</b>	<b>10</b>

<i>Cube Aggregation</i>	10
<i>Cube Operators</i>	11
<i>Extended Cube</i>	11
<i>Aggregation Using Hierarchies</i>	11
<i>Pivoting</i>	12
<i>CUBE Operator (SQL-99)</i>	12
<i>Aggregates</i>	12
<i>Query &amp; Analysis Tools</i>	12
<i>Other Operations</i>	13
<i>Data Warehouse Implementation</i>	13
<i>Multi-Tiered Architecture</i>	13
<i>Monitoring</i>	13
<i>Data Cleansing</i>	14
<i>Loading data</i>	14
<i>OLAP Implementation</i>	14
<i>Derived Data</i>	14
<i>What to Materialize</i>	14
<i>Materialization Factors</i>	14
<i>Cube Aggregates Lattice</i>	15
<i>Dimension Hierarchies</i>	15
<i>Interesting Hierarchy</i>	16
<i>Indexing OLAP Data: Bitmap Index</i>	16
<i>Join Processing</i>	16
<i>Schema for Examples</i>	16
<i>Equality Joins with One Join Column</i>	17
<i>Simple Nested Loops Join</i>	17
<i>Block Nested Loops Join</i>	17
<i>Examples of Block Nested Loops</i>	18
<i>Index Nested Loops Join</i>	18
<i>Examples of Index Nested Loops</i>	18
<i>Sort-Merge Join</i>	19
<i>Example of Sort-Merge Join</i>	19
<i>Refinement of Sort-Merger Join</i>	19
<i>Hash-Join</i>	20
<i>Observations in Hash-Join</i>	20
<i>Cost of Hash-Join</i>	20

<i>Join Indices</i> .....	20
<i>General Join Conditions</i> .....	21
<i>An invention in 2000s: Column Stores for OLAP</i> .....	21
<i>Row Store and Column Store</i> .....	21
<i>Column Stores: High Level</i> .....	22
<i>Data Model (Vertica/C-Store)</i> .....	22
<i>Read store: Column Encoding/Compression</i> .....	23
<i>Compression</i> .....	24
<i>Write Store</i> .....	24
<i>How to solve read/write conflict</i> .....	24
<i>Query Execution – Operators</i> .....	24
<i>Benefits in Query Processing</i> .....	24
<i>Evaluation</i> .....	24
<i>Query Performance</i> .....	25
<i>Summary: The performance gain</i> .....	25
<i>Google's Dremel/Big Query: Interactive Analysis of Web-Scale Datasets</i> .....	25
<i>Big Query System</i> .....	25
<i>Widely Used inside Google</i> .....	26
<i>Nested Data Model</i> .....	27
<i>Column-striped representation</i> .....	27
<i>Repetition and definition levels</i> .....	28
<i>Query processing</i> .....	28
<i>SQL dialect for nested data</i> .....	28
<i>Serving tree</i> .....	29
<i>Example: count()</i> .....	29
<i>Experiments</i> .....	29
<i>Interactive Speed</i> .....	30
<i>Big Query: Powered by Dremel</i> .....	30
<i>List of Column Data Bases</i> .....	30
<i>Take-home messages</i> .....	30



1



2

---

# Agenda

- History of Spark
- Introduction
- Components of Stack
- Resilient Distributed Dataset – RDD



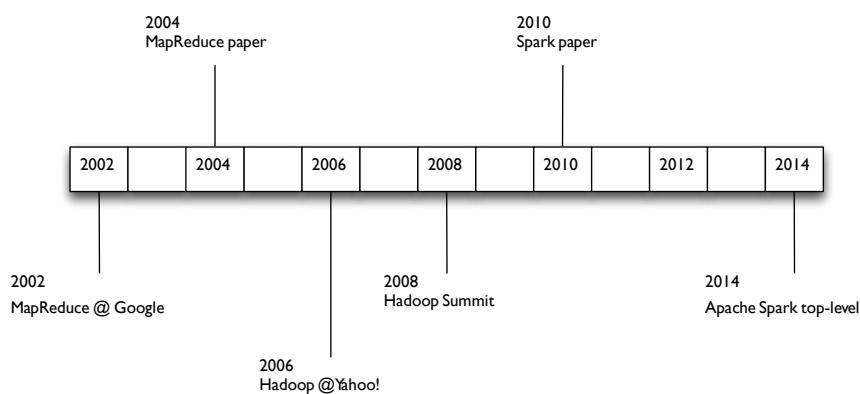
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

3

---

---

## History of Spark



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

4

# History of Spark

circa 1979 – **Stanford, MIT, CMU**, etc.

set/list operations in LISP, Prolog, etc., for parallel processing  
[www-formal.stanford.edu/jmc/history/lisp/lisp.htm](http://www-formal.stanford.edu/jmc/history/lisp/lisp.htm)

circa 2004 – **Google**

*MapReduce: Simplified Data Processing on Large Clusters* Jeffrey Dean and Sanjay Ghemawat  
[research.google.com/archive/mapreduce.html](http://research.google.com/archive/mapreduce.html)

circa 2006 – **Apache**

*Hadoop*, originating from the Nutch Project Doug Cutting  
[research.yahoo.com/files/cutting.pdf](http://research.yahoo.com/files/cutting.pdf)

circa 2008 – **Yahoo**

*web scale search indexing Hadoop Submit, HUG, etc.*  
[developer.yahoo.com/hadoop/](http://developer.yahoo.com/hadoop/)

circa 2009 – **Amazon AWS**

Elastic MapReduce  
Hadoop modified for EC2/S3, plus support for Hive, Pig, Cascading, etc.  
[aws.amazon.com/elasticmapreduce/](http://aws.amazon.com/elasticmapreduce/)

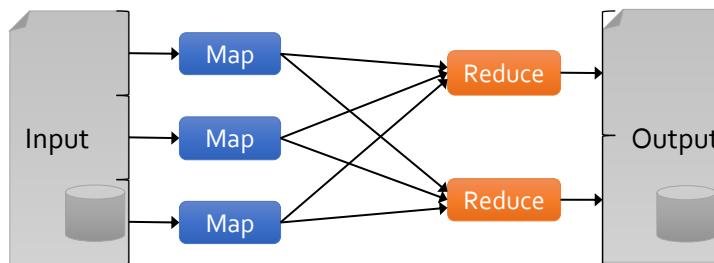


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

5

# MapReduce

Most current cluster programming models are based on *acyclic data flow* from stable storage to stable storage



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

6

# MapReduce

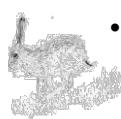
- Acyclic data flow is inefficient for applications that repeatedly reuse a *working set* of data:
  - **Iterative** algorithms (machine learning, graphs)
  - **Interactive** data mining tools (R, Excel, Python)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

7

## Data Processing Improvement Goals



- **Low latency (interactive) queries on historical data:** enable faster decisions
  - E.g., identify why a site is slow and fix it
- **Low latency queries on live data (streaming):** enable decisions on real-time data
  - E.g., detect & block worms in real-time (a worm may infect 1mil hosts in 1.3sec)
- **Sophisticated data processing:** enable “better” decisions
  - E.g., anomaly detection, trend analysis



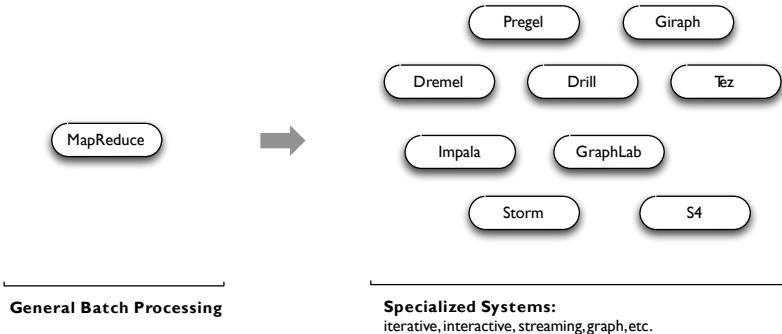
Therefore, people built specialized  
systems as workarounds...



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

8

# Specialized Systems



*The State of Spark, and Where We're Going Next*

**Matei Zaharia**

Spark Summit (2013)

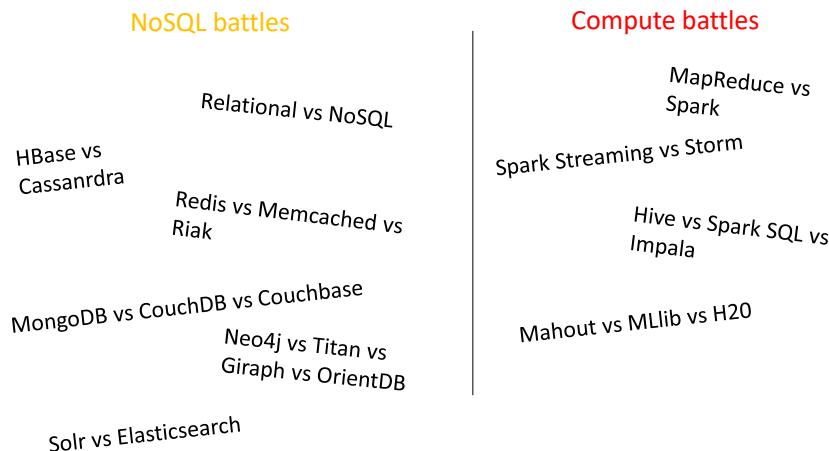
[youtu.be/nU6vO2EJAb4](https://youtu.be/nU6vO2EJAb4)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

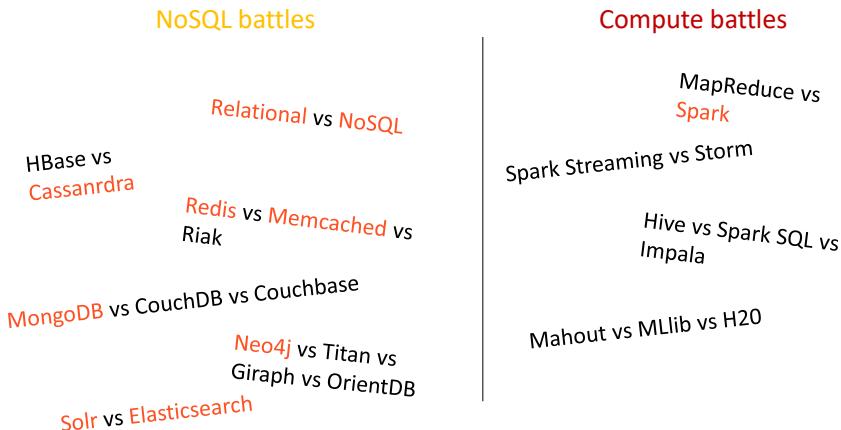
9

# Storage vs Processing Wars



10

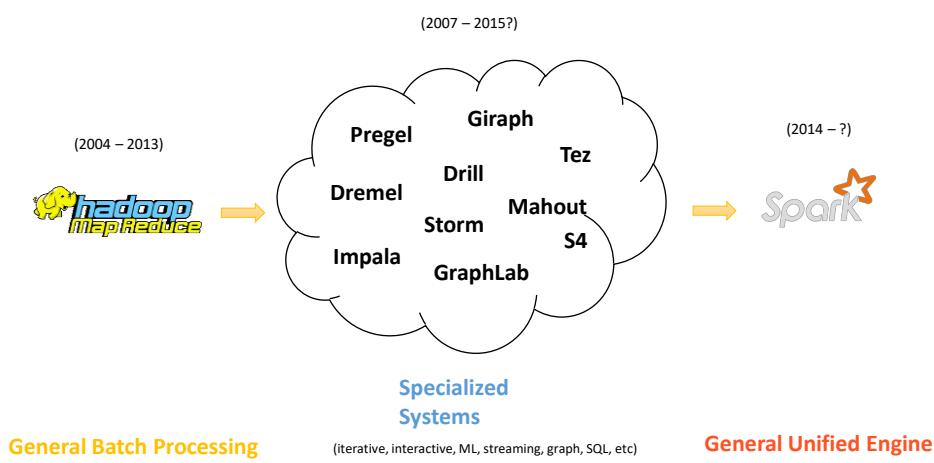
# Storage vs Processing Wars



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

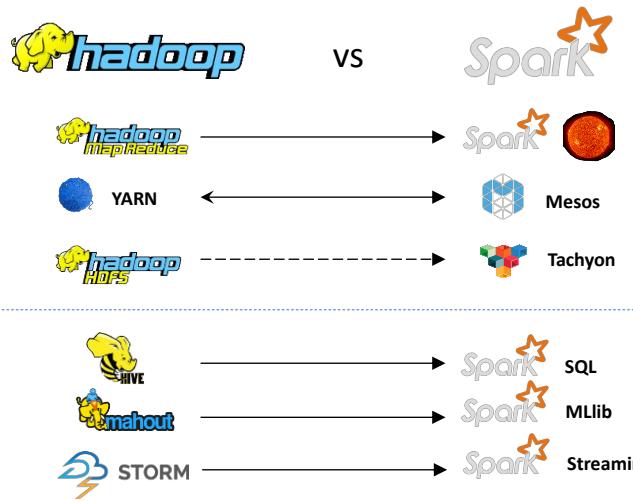
11

# Specialized Systems



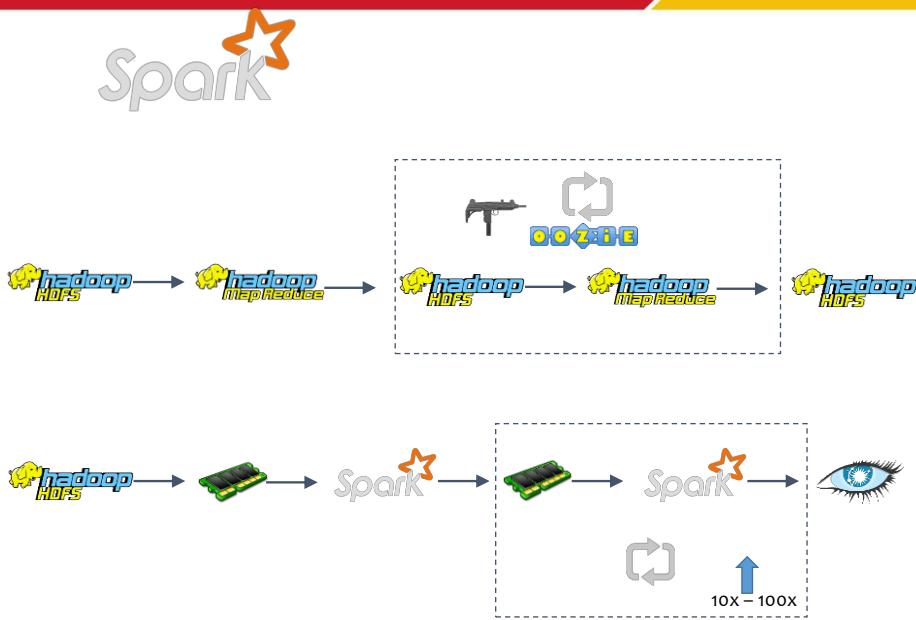
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

12



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

13

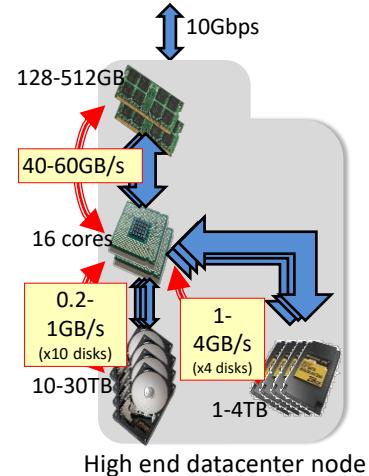


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

14

# Support Interactive and Streaming Comp.

- Aggressive use of ***memory***
- Why?
  1. Memory transfer rates >> disk or SSDs
  2. Many datasets already fit into memory
    - Inputs of over 90% of jobs in Facebook, Yahoo!, and Bing clusters fit into memory
    - e.g., 1TB = 1 billion records @ 1KB each
  3. Memory density (still) grows with Moore's law
    - RAM/SSD hybrid memories at horizon

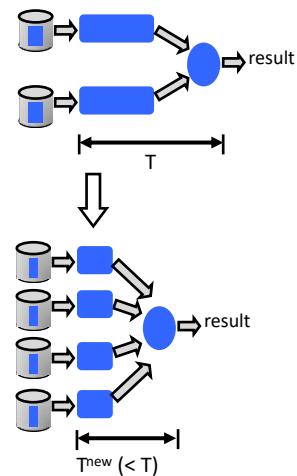


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

15

# Support Interactive and Streaming Comp.

- Increase ***parallelism***
- Why?
  - Reduce work per node → improve latency
- Techniques:
  - Low latency parallel **scheduler** that achieve high locality
  - Optimized **parallel communication patterns** (e.g., shuffle, broadcast)
  - Efficient **recovery** from failures and straggler mitigation

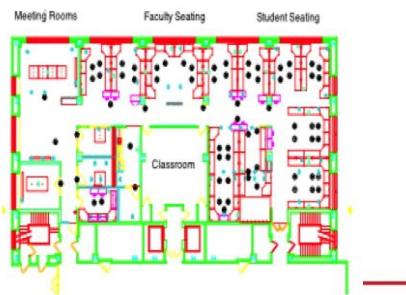


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

16

# Berkeley AMPLab

- “Launched” January 2011: 6 Year Plan
- 8 CS Faculty
- ~40 students
- 3 software engineers
- Organized for collaboration:



17

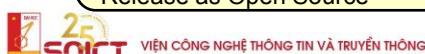
# Berkeley AMPLab

- Funding:
    - DABPA, KData, NSF CISE Expedition Grant
    - Industrial, founding sponsors
    - 18 other sponsors, including
- 

**Goal:** Next Generation of Analytics Data Stack for Industry &

Research:

- Berkeley Data Analytics Stack (BDAS)
- Release as Open Source

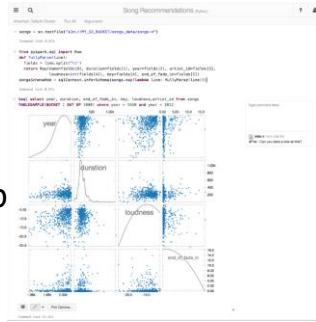


18

# Databricks



- Founded in late 2013
- by the creators of Apache Spark
- Original team from UC Berkeley AMPLab
- Raised \$47 Million in 2 rounds



Databricks Cloud:  
"A unified platform for building Big Data pipelines – from ETL to Exploration and Dashboards, to Advanced Analytics and Data Products."



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

19

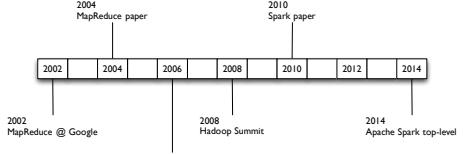
The Databricks team contributed more than **75%** of the code added to Spark in the 2014



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

20

# History of Spark



**Spark: Cluster Computing with Working Sets**  
Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, Ion Stoica  
University of California, Berkeley

**Abstract**  
MapReduce and its variants have been highly successful in implementing large-scale data-processing applications in distributed systems. However, these frameworks are built around an acyclic data flow model that is unsuitable for many applications. In this paper, we propose a new abstraction for working sets that allows us to express iterative workloads on such class of applications: those that make multiple passes over the same dataset to refine their results. This includes many iterative machine learning algorithms, as well as many other applications such as graph processing. Our new framework called Spark that supports these applications achieves significant performance improvements over MapReduce. To achieve these goals, Spark introduces an abstraction called RDD (Resilient Distributed Dataset), which is a fault-tolerant collection of objects partitioned across a cluster of machines. An RDD can be thought of as a generalization of a dataset in that it can be updated incrementally. Spark can execute programs by its iterative nature, and can even support complex queries with sub-second latency. It can also support distributed machine learning operations with sub-second response time.

**Introduction**  
A new model of cluster computing has become widely popular, in which data-parallel computations are executed on clusters of commodity machines. These systems automatically provide load-balancing, scheduling, fault tolerance, and distributed memory management. Two prominent examples of this model are MapReduce [17] and MapReduce-like systems such as Dryad [7] and Hadoop [2]. These systems achieve their scalability and fault tolerance by providing a simple interface for iterative applications. An RDD is a fault-tolerant collection of objects partitioned across a cluster of machines. An RDD can be thought of as a generalization of a dataset in that it can be updated incrementally. Spark can execute programs by its iterative nature, and can even support complex queries with sub-second latency. It can also support distributed machine learning operations with sub-second response time.

**Background**  
Many common machine learning algorithms require iterative passes over the same dataset to optimize a parameter (e.g., through gradient descent). While each iteration can be expressed as a

MapReduce [17] is just one of many distributed systems that have achieved significant performance improvements over MapReduce. Many others, such as Hadoop [2], have followed a similar path of evolution, while providing similar capabilities and fault tolerance.

**Iterative analytics:** Hadoop is often used to run iterative analytics. For example, Pig [21] and Hive [12] build upon Hadoop to support iterative data processing. These frameworks work on a number of machines and query it sequentially, which leads to significant latency times of seconds because it runs an iterative algorithm sequentially.

This paper presents a new cluster computing framework called Spark that provides a fault-tolerant abstraction for iterative workloads while providing similar capabilities and fault tolerance.

The main abstraction in Spark is that of a **RDD**, a fault-tolerant distributed dataset. An RDD is a collection of objects partitioned across a cluster of machines, such as a HDFS file or a database. An RDD can be thought of as a generalization of a dataset in that it can be updated incrementally. Spark is implemented in Scala [21], a mostly typed functional programming language. Scala is a superset of Java and allows for mixed-mode development. Scala supports a functional programming interface similar to Scheme and a general-purpose object-oriented interface similar to C++, and allows for mixed-mode development. Although our implementation of Spark is still in its early stages, we believe that it will be competitive with existing distributed systems. Although our implementation of Spark is still in its early stages, we believe that it will be competitive with existing distributed systems.

This paper is organized as follows. Section 2 describes how to use RDDs to implement iterative applications. We then present a detailed description of the RDD abstraction in Section 3. Section 4 discusses the implementation of Spark. Section 5 concludes this paper.

## Spark: Cluster Computing with Working Sets

Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, Ion Stoica  
USENIX HotCloud (2010)

[people.csail.mit.edu/matei/papers/2010/hotcloud\\_spark.pdf](http://people.csail.mit.edu/matei/papers/2010/hotcloud_spark.pdf)

## Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing

Matei Zaharia, Mosharaf Chowdhury, Taghagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica NSDI (2012)

[usenix.org/system/files/conference/nsdi12/nsdi12-final138.pdf](http://usenix.org/system/files/conference/nsdi12/nsdi12-final138.pdf)

21



# History of Spark

**Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing**  
Matei Zaharia, Mosharaf Chowdhury, Taghagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica  
University of California, Berkeley

**Abstract**  
We present Resilient Distributed Datasets (RDDs), a distributed memory abstraction that lets programmers handle large-scale data-processing tasks in a fault-tolerant manner. RDDs are motivated by two types of applications: iterative algorithms and interactive data mining tools. Iterative algorithms are common in machine learning, and can improve performance by an order of magnitude by reusing data in memory. Interactive data mining tools are common in real-time systems, and benefit from a restricted form of shared memory, based on coarse-grained transformations rather than fine-grained updates. These systems are more likely to reuse data than to read it from disk. In this paper, we propose a new abstraction called RDD, which provides a fault-tolerant abstraction for in-memory storage. RDDs are designed to support both iterative and interactive data reuse in a broad range of applications. RDDs are fault-tolerant, and can be updated incrementally. They let programmers explicitly persist intermediate results in memory, control their partitioning to optimize data placement, and map them to parallel execution units.

**Introduction**  
Cluster computing frameworks like MapReduce [10] and Dryad [7] have been very successful in implementing large-scale data-processing applications. These systems are built using parallel computation primitives such as map and reduce, and are designed to support a wide variety of computations, including iterative algorithms, distributed machine learning jobs, and interactive data mining tools. Although these frameworks provide numerous advantages, they lack abstraction for leveraging distributed memory. They also lack abstraction for dealing with the class of computing applications that reuse intermediate data results. For example, iterative machine learning algorithms such as many iterative machine learning and graph algorithms, and interactive data mining tools like linear and logistic regression. Another use case is interactive data mining tools, such as data warehouses that query the same subset of the data. Definition: In both cases, keeping data in memory can improve performance by an order of magnitude.”

Recognizing this problem, researchers have developed various distributed memory abstractions for in-memory data reuse. For example, Pig [22] is a system for iterative data mining that stores intermediate results in memory, while Hadoop [7] offers a memory MapReduce interface. However, these frameworks only support a limited set of applications. In contrast, RDDs support a wide variety of applications, including iterative MapReduce steps, and perform data sharing implicitly via the RDD abstraction. In addition, RDDs support more general modes, e.g., to fit a wide variety of datasets into memory, and to support different access patterns.

In this paper, we propose a new abstraction called RDD, which provides a fault-tolerant abstraction for in-memory storage. RDDs are designed to support both iterative and interactive data reuse in a broad range of applications. RDDs are fault-tolerant, and can be updated incrementally. They let programmers explicitly persist intermediate results in memory, control their partitioning to optimize data placement, and map them to parallel execution units.

The main challenge in designing RDDs is how to provide fault tolerance in a fault-tolerant way. This requires a trade-off between fault tolerance and fault tolerance efficiency. Existing systems provide an iterative model for fault tolerance, where the only way to recover from a failure is to recompute the entire dataset from scratch. In contrast, these systems RDDs provide an interactive model for fault tolerance, where the only way to recover from a failure is to log the transformations used to build a dataset, and then to replay them on a new machine. This is particularly important for distributed memory applications, since checkpointing the data in memory may be much slower than page-churn growth, however, and we discuss how to do it in §5.

**“We present Resilient Distributed Datasets (RDDs), a distributed memory abstraction that lets programmers perform in-memory computations on large clusters in a fault-tolerant manner.**

**RDDs are motivated by two types of applications that current computing frameworks handle inefficiently: iterative algorithms and interactive data mining tools.**

**In both cases, keeping data in memory can improve performance by an order of magnitude.”**



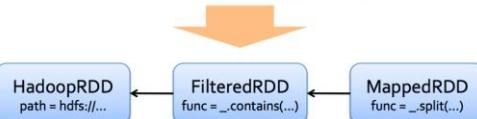
22

# History of Spark

## RDD Fault Tolerance

RDDs track the series of transformations used to build them (their *lineage*) to recompute lost data

E.g: `messages = textFile(...).filter(_.contains("error")).map(_.split('\t')(2))`



*The State of Spark, and Where We're Going Next*

Matei Zaharia

Spark Summit (2013)



23

# History of Spark



Analyze real time streams of data in  $\frac{1}{2}$  second intervals

[www.cs.berkeley.edu/~matei/papers/2013/soosp\\_spark\\_streaming.pdf](http://www.cs.berkeley.edu/~matei/papers/2013/soosp_spark_streaming.pdf)

**Discretized Streams: Fault-Tolerant Streaming Computation at Scale**  
Matei Zaharia, Tathagata Das, Haoyuan Li, Timothy Hunter, Scott Shenker, Ion Stoica  
University of California, Berkeley

**Abstract**  
Many “big data” applications must act on data in real time. Running these applications at ever-larger scales requires parallel platforms that automatically handle faults and mitigate performance variability. In this paper, we propose a fault-tolerant streaming computation model that provides fault recovery in an inexpensive manner, requiring hot replication or long recovery times only for a small fraction of the system. Our processing model, *discretized streams* (DStreams), thus overcomes these challenges. DStreams enable a parallel real-time processing of distributed data streams using traditional replication and bucking schemes, and tolerates stragglers. We show that they support a rich set of operators, where each operator can fail independently and scale to single-node systems, linear scaling to 100 nodes, sub-second recovery from failure, and fast fault recovery. DStreams can easily be composed with batch and interactive query models like MapReduce, enabling rich applications that span all these modes. We implement DStreams in a system called Spark Streaming.

**1 Introduction**  
Much of “big data” is received in real time, and is most valuable at its time of arrival. For example, a social network may wish to detect trending conversation topics in

faults and anomalies (like nodes). Both problems are inevitable in large clusters [12], so streaming applications must recover from them quickly. Fast recovery is even more important for real-time data, as it was shown that a 30-second delay to recover from a fault or straggler is a misfortune in a batch setting, if it meant being delayed for a day [2].

Unfortunately, existing streaming systems have limited fault and straggler tolerance. Most distributed systems, such as Storm [17], Flink [27], and Twitter Stream [33], MapReduce Online [11], and streaming databases [5, 9, 10], are based on a continuous operator model, where each node receives data, processes it, receive each record, update internal state, and send new records. While this model is quite natural, it makes it difficult to tolerate faults and stragglers.

Specifically, given the continuous operator model, systems perform many more tasks than necessary [20], especially when there are two copies of each node [5, 34], or *apartition backup*, where nodes buffer sent messages and periodically flush them to disk [33, 11, 37]. Neither approach is attractive in large clusters: replication costs 2× the hardware, while sparsify backup takes a long time to recover, as the whole system must wait for a new node to serially rebuild the partition.

```

TwitterUtils.createStream(...)  
.filter(_.getText.contains("Spark"))  
.countByWindow(Seconds(5))
  
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

24

# History of Spark



Seemlessly mix SQL queries with Spark programs.

## Spark SQL: Relational Data Processing in Spark

Michael Armbrust<sup>1</sup>, Reynold S. Xin<sup>1</sup>, Cheng Lian<sup>1</sup>, Yin Hua<sup>1</sup>, Davies Liu<sup>1</sup>, Joseph K. Bradley<sup>1</sup>, Xiangru Meng<sup>1</sup>, Tomer Katta<sup>1</sup>, Michael J. Franklin<sup>1</sup>, Ali Ghodsi<sup>2</sup>, Matei Zaharia<sup>1</sup>  
<sup>1</sup>Databricks Inc.    <sup>2</sup>AMPLab, UC Berkeley

### ABSTRACT

Spark SQL is a new module in Apache Spark that integrates processing of structured data with the core API. Built on our experience with Struct, Spark SQL lets Spark programs benefit from the strengths of both SQL and Spark. It supports queries and optimized storage, and let SQL users write full complex analytic programs in SQL. Unlike previous systems, Spark SQL makes two main additions. First, it offers native support for relational data processing and storage, processing through a declarative DataFrame API that integrates with preexisting Spark code. Second, it includes a highly extensible optimizer that translates SQL queries into Spark programs, a language that makes it easy to add comparable roles, control code generation, and reuse existing code. We demonstrate that Spark SQL builds a variety of features (e.g., schema inference for JSON, machine learning and graph processing) on top of the core DataFrame API for the complex needs of modern data analysis. We see Spark SQL as an evolution of both SQL-on-Spark and of Spark itself, offering take advantage of both while leveraging the benefits of the Spark programming model.

### Categories and Subject Descriptors

H.2 [Database Management]: Systems

### Keywords

Databases; Data Warehouses; Machine Learning; Spark; Hadoop

### 1 Introduction

Big data applications require a mix of processing techniques, data sources and storage formats. The earliest systems designed for these workloads, such as MapReduce, were sites a powerful but

While the popularity of relational systems shows that more often prefer writing declarative operators, the relational approach is insufficient for many big data applications. First, users want to perform ETL (extract, transform, load) operations on data that is not well-structured, requiring custom code. Second, users want to perform advanced analytics, such as machine learning and graph processing, that are better suited to iterative computation. Finally, we have observed that more data pipelines would simply be easier to maintain if they were built on top of a declarative procedural algorithm. Unfortunately, these two classes of systems—declarative and procedural—have traditionally been largely forcing users to choose one paradigm or the other.

This paper describes our effort to combine both worlds in Spark SQL, a new module in the core Spark API. Built on our experience with Struct, Spark SQL lets users seamlessly interact with the core DataFrame API.

Spark SQL uses the parallel processing model from the DataFrame API to perform relational operations on both external data sources and in-memory datasets. By supporting both the declarative style of the widely used data frame concept in R [13], and evaluates operations locally, Spark SQL provides a familiar interface for users to support the wide range of data sources and algorithms in big data. Spark SQL also provides a rich set of optimization rules, such as Catalyst makes it easy to add data sources, optimization rules, and data types for domains such as machine learning.

The paper is organized as follows. Section 2 gives a brief procedural integration within Spark programs. Databases are collections of structured metadata that can be queried and updated. In previous sections, we used new introduced APIs that allow better optimizations. They can

```
sqlCtx = new HiveContext(sc)
results = sqlCtx.sql(
    "SELECT * FROM people")
names = results.map(lambda p:
    p.name)
```

25

# History of Spark



Analyze networks of nodes and edges using graph processing

## GraphX: A Resilient Distributed Graph System on Spark

Reynold S. Xin, Joseph E. Gonzalez, Michael J. Franklin, Ion Stoica  
 AMPLab, EECS, UC Berkeley  
 {rxin, jgonzal, franklin, ionstoic}@cs.berkeley.edu

### ABSTRACT

From social networks to surgical advertising, big graphs capture the essence of data and are central to recent advances in machine learning and data mining. Unfortunately, directly applying existing distributed systems to graphs is inefficient. The need for iterative, scalable tools for graph computation has led to the development of new graph-parallel systems (e.g., Giraph, GraphX). These systems are designed to efficiently execute graph algorithms. Unfortunately, these new graph-parallel systems are often just as problematic as the subsequent computational primitives they are built upon. Specifically, they lack fault tolerance and support for interactive data mining.

We introduce GraphX, a distributed system for distributed graph-parallel and graph-parallel systems by efficiently expressing graph programs. Specifically, GraphX provides a framework to merge new ideas in distributed graph processing with efficiently distribute graphs as稀疏 data-structures. Similarly, we leverage advances in distributed systems to build a system that is both fast and fault-tolerant. We provide powerful new operations to simplify graph computation. Finally, we demonstrate how GraphX can implement the PowerGraph and Pregel abstractions in less than 20 lines of code. By providing a clean interface to these primitives, we enable users to interactively load, transform, and compute on massive graphs.

### 1 INTRODUCTION

From social networks to advertising and the web, big graphs can

be found in a wide range of important applications. By modeling the

and distributed systems. By abstracting away the challenges of large-scale distributed systems, we hope that this paper will inspire further development, implementation, and application of sophisticated graph algorithms to large-scale real-world graph problems.

Graphs are a natural way to represent many common properties, such as proximity and connections. As a result, graph computation is a fundamental building block for many common graph algorithms and applications. Unfortunately, because each node in a graph is connected to many other nodes, these challenges are often compounded by the inherent nature of graph computation.

Furthermore, while these frameworks address

the challenges of graph computation, they do not address the challenges of interpreting and applying the results of computation. Finally, few frameworks support distributed data structures.

Alternatively, distributed systems like MapReduce and Spark have been developed to support iterative computation suited to the task of graph computation (ETL). By exploiting the strengths of these distributed systems, we can build a system that supports a range of fault-tolerance strategies.

More recent systems like GraphX have attempted to merge distributed data structures with distributed data flows. However, merging distributed data structures with distributed data flows can be challenging and typically leads to significant performance overheads.

Finally, another challenge we introduce GraphX, a graph-computation system which runs on the Spark distributed framework.

GraphX provides a distributed graph abstraction and distributed structure to introduce the Resilient Distributed Graph (RDG), which associates records with vertices and edges in a graph and provides a collection of expressive computational primitives. Using these

```
graph = Graph(vertices, edges)
messages =
spark.textFile("hdfs://...")
graph2 =
graph.joinVertices(messages) {
  (id, vertex, msg) => ...
}
```

[https://amplab.cs.berkeley.edu/wp-content/uploads/2013/05/grades-graphx\\_with\\_fonts.pdf](https://amplab.cs.berkeley.edu/wp-content/uploads/2013/05/grades-graphx_with_fonts.pdf)

26



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

# History of Spark



SQL queries with Bounded Errors and Bounded Response Times

## BlinkDB: Queries with Bounded Errors and Bounded Response Times on Very Large Data

Sameer Agarwal<sup>1</sup>, Barzan Mozafari<sup>2</sup>, Aurojit Panda<sup>1</sup>, Henry Müller<sup>1</sup>, Samuel Madden<sup>1</sup>, Jon Stoica<sup>1\*</sup>  
<sup>1</sup>University of California, Berkeley    <sup>2</sup>Massachusetts Institute of Technology    \*Conviva Inc.  
{sameerag, apanda, henrym, istoica}@cs.berkeley.edu, {barzan, madden}@csail.mit.edu

### Abstract

In this paper, we present BlinkDB, a massively parallel, approximate query engine for running interactive SQL queries on large datasets. BlinkDB allows users to trade off query accuracy for response time, allowing interactive queries over massive data by running queries on data samples and presenting results annotated with meaningful error bars. To achieve this, BlinkDB consists of two main parts: (1) an adaptive optimization framework that builds and refines a set of multi-dimensional stratified samples from original data over time, and (2) a dynamic sample selection strategy that selects appropriate samples to answer a query given its response time requirements. We evaluate BlinkDB against the well-known TPC-H benchmarks and a real-world analytic workload derived from Conviva Inc., a company that manages video delivery networks. Our experiments on a 100 node cluster show that BlinkDB can answer queries on up to 27 TBs of data in less than 2 seconds (over 200x faster than Hive), within an error of >10%.

### 1. Introduction

Modern data analytics applications involve computing aggregates over a large number of records to roll-up web clicks,

co-occurring of large amounts of data by trading result accuracy for response time and space. These techniques include sampling [10, 14], sketches [12] and on-line aggregation [1]. To understand how these techniques work, consider the following simple query that computes the average `SessionTime` over all users originating in New York:

`SELECT AVG(SessionTime)`

`FROM Sessions`

`WHERE SessionCity = "New York"`

Suppose the `Sessions` table contains 100 million tuples for New York, and cannot fit in memory. In that case, the above query may take a long time to execute, since disk reads are expensive, and such a query needs multiple disk accesses to read all the tuples. Suppose instead that we have executed the same query on a sample containing only 10,000 New York tuples, such that the entire sample fits in memory. This will be orders of magnitude faster, while the process of aggregation will incur a few percent of the actual value, an accuracy good enough for many practical purposes. Using sampling theory we could even provide confidence bounds on the accuracy of the answer [6].

Precisely described approximation techniques make different trade-offs between efficiency and the generality of the

```
SELECT avg(sessionTime)
FROM Table
WHERE city='San Francisco'
WITHIN 2 SECONDS
```

Queries with Time Bounds

```
SELECT avg(sessionTime)
FROM Table
WHERE city='San Francisco'
ERROR 0.1 CONFIDENCE 95.0%
```

Queries with Error Bounds

[https://www.cs.berkeley.edu/~sameerag/blinkdb\\_eurosys13.pdf](https://www.cs.berkeley.edu/~sameerag/blinkdb_eurosys13.pdf)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

27

# History of Spark

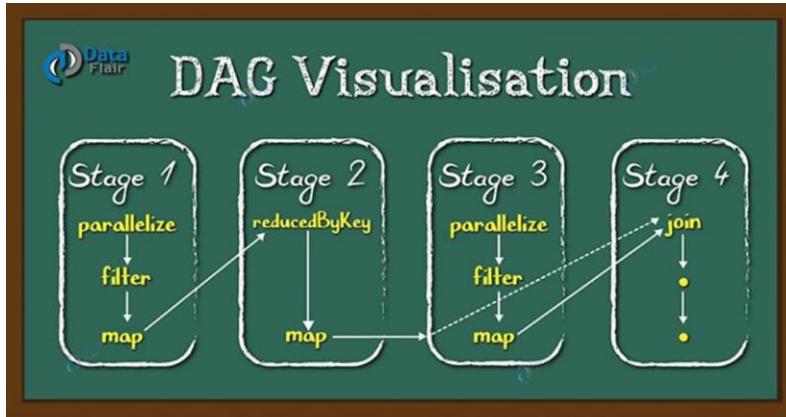
- Unlike the various specialized systems, Spark's goal was to *generalize* MapReduce to support new apps within same engine
- Two reasonably small additions are enough to express the previous models:
  - *fast data sharing*
  - *general DAGs*
- This allows for an approach which is more efficient for the engine, and much simpler for the end users



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

28

## Directed Acyclic Graph - DAG

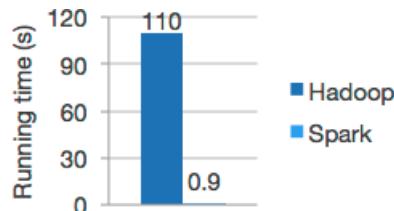


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

29

## What is Apache Spark

- Spark is a unified analytics engine for large-scale data processing
- Speed:** run workloads 100x faster
  - High performance for both batch and streaming data
  - Computations run in memory



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

30

# What is Apache Spark

- **Ease of Use:** write applications quickly in Java, Scala, Python, R, SQL
  - Offer over 80 high-level operators
  - Use them interactively from Scala, Python, R, and SQL

```
df = spark.read.json("logs.json")
df.where("age > 21")
select("name.first").show()
```

Spark's Python DataFrame API  
Read JSON files with automatic schema inference

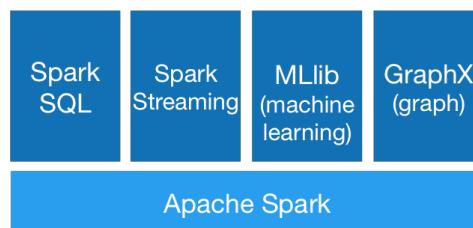


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

31

# What is Apache Spark

- **Generality:** combine SQL, Streaming, and complex analytics
  - Provide libraries including SQL and DataFrames, Spark Streaming, MLlib, GraphX
  - Wide range of workloads e.g., batch applications, interactive algorithms, interactive queries, streaming



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

32

# What is Apache Spark

- Run Everywhere:

- run on Hadoop, Apache Mesos, Kubernetes, standalone or in the cloud.
- access data in HDFS, Aluxio, Apache Cassandra, Apache Hbase, Apache Hive, etc.



## Comparison between Hadoop and Spark

	Hadoop	Spark
<b>Strengths</b>	<ul style="list-style-type: none"> <li>Can collect any data</li> <li>Limitless in size</li> </ul>	<ul style="list-style-type: none"> <li>Can work off any Hadoop collection</li> <li>Runs on Hadoop, or other clusters</li> <li>In-memory processing makes it very fast</li> <li>Supports Java, Scala, Python, and R*, and can be used with SQL.</li> </ul>
<b>Used for</b>	<ul style="list-style-type: none"> <li>Initial data ingestion</li> <li>Data curation</li> <li>Large-scale “boil the ocean” analytics</li> <li>Data archiving</li> </ul>	<ul style="list-style-type: none"> <li>Complex query processing of large amounts of data quickly</li> <li>Can handle ad hoc queries</li> </ul>
<b>Limitations</b>	<ul style="list-style-type: none"> <li>MapReduce is hard to program</li> <li>Disk-based batch nature limits speed, agility.</li> </ul>	<ul style="list-style-type: none"> <li>Limited only by processor speed, available memory, cores, and cluster size.</li> </ul>

# 100TB Daytona Sort Competition

	Hadoop MR Record	Spark Record	Spark 1 PB
Data Size	102.5 TB	100 TB	1000 TB
Elapsed Time	72 mins	23 mins	234 mins
# Nodes	2100	206	190
# Cores	50400 physical	6592 virtualized	6080 virtualized
Cluster disk throughput	3150 GB/s (est.)	618 GB/s	570 GB/s
Sort Benchmark Daytona Rules	Yes	Yes	No
Network	dedicated data center, 10Gbps	virtualized (EC2) 10Gbps network	10Gbps network
Sort rate	1.42 TB/min	4.27 TB/min	4.27 TB/min
Sort rate/node	0.67 GB/min	20.7 GB/min	22.5 GB/min

Spark sorted the same data **3X faster** using **10X fewer machines** than Hadoop MR in 2013.

All the sorting took place on disk (HDFS)  
without using Spark's in-memory cache!



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

35

WIRED

GEAR SCIENCE ENTERTAINMENT BUSINESS SECURITY DESIGN OPINION MAGAZINE

ENTERPRISE

big data · databricks · google · Hadoop

## Startup Crunches 100 Terabytes of Data in a Record 23 Minutes

BY KIINT FINLEY 10.13.14 | 2:36 PM | PERMALINK

[Share](#) 1.1k [Tweet](#) 789 [G+1](#) 75 [in Share](#) 565 [Pin it](#)

**GIGAOM** EVENTS RESEARCH

Cloud Data Media Mobile Science & Energy Social & Web Podcasts

Gigaom Research. Get unlimited market intelligence from over 200+

MUST READS

Google launches Contributor, a crowdfunding tool for publishers

Net neutrality looks doomed in Europe before it even gets started

Five tech products that designers have fallen in love with

Databricks demolishes big data benchmark to prove Spark is fast on disk, too

By Derrick Harris Oct 10, 2014 - 149 PM PST

[Comment](#)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

36

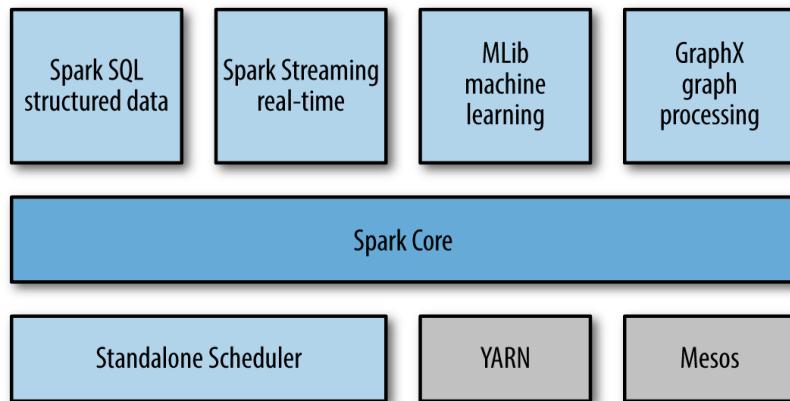
# Components of Stack



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

37

## The Spark stack



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

38

19

# The Spark stack

- Spark Core:
  - contain basic functionality of Spark including task scheduling, memory management, fault recovery, etc.
  - provide APIs for building and manipulating RDDs
- SparkSQL
  - allow querying structured data via SQL, Hive Query Language
  - allow combining SQL queries and data manipulations in Python, Java, Scala



39

# The Spark stack

- Spark Streaming: enables processing of live streams of data via APIs
- Mlib:
  - contain common machine language functionality
  - provide multiple types of algorithms: classification, regression, clustering, etc.
- GraphX:
  - library for manipulating graphs and performing graph-parallel computations
  - extend Spark RDD API



40

# The Spark stack

- Cluster Managers
  - Hadoop Yarn
  - Apache Mesos, and
  - Standalone Scheduler (simple manager in Spark).



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

41

## Resilient Distributed Dataset – RDD

- RDD Basics
- Creating RDDs
- RDD Operations
- Common Transformation and Actions
- Persistence (Caching)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

42

# RDD Basics

- RDD:
  - Immutable distributed collection of objects
  - Split into multiple partitions => can be computed on different nodes
- All work in Spark is expressed as
  - creating new RDDs
  - transforming existing RDDs
  - calling actions on RDDs

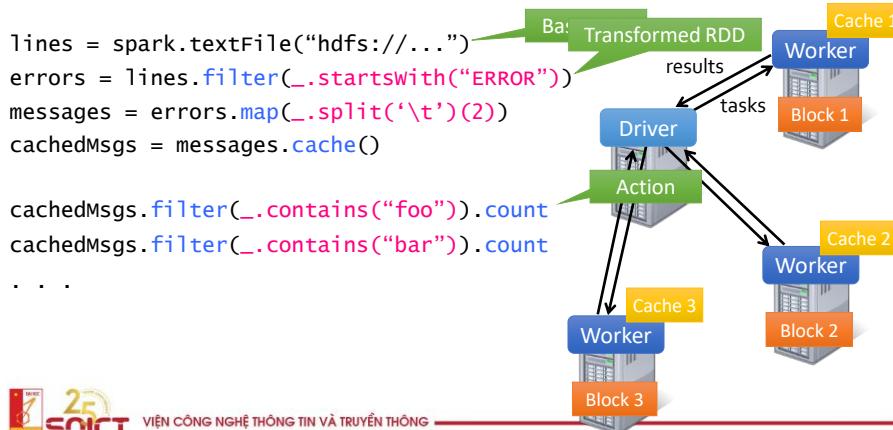


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

43

## Example

Load error messages from a log into memory, then interactively search for various patterns

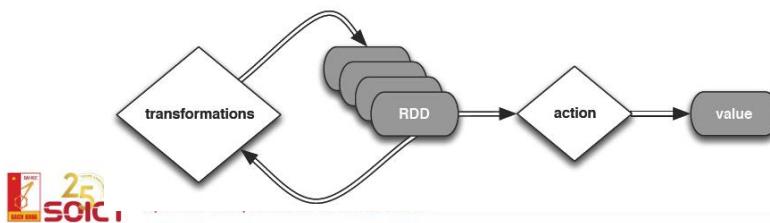


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

44

# RDD Basics

- Two types of operations: transformations and actions
- Transformations: construct a new RDD from a previous one e.g., filter data
- Actions: compute a result base on an RDD e.g., count elements, get first element



45

## Transformations

- Create new RDDs from existing RDDs
- Lazy evaluation
  - See the whole chain of transformations
  - Compute just the data needed
- Persist contents:
  - persist an RDD in memory to reuse it in future
  - persist RDDs on disk is possible

# Typical works of a Spark program

1. Create some input RDDs from external data
2. Transform them to define new RDDs using transformations like filter()
3. Ask Spark to persist() any intermediate RDDs that will need to be reused
4. Launch actions such as count(), first() to kick off a parallel computation



47

## Resilient Distributed Dataset – RDD

- RDD Basics
- Creating RDDs
- RDD Operations
- Common Transformation and Actions
- Persistence (Caching)



48

## Two ways to create RDDs

### 1. Parallelizing a collection: uses parallelize()

- Python

```
lines = sc.parallelize(["pandas", "i like pandas"])
```

- Scala

```
val lines = sc.parallelize(List("pandas", "i like
pandas"))
```

- Java

```
JavaRDD<String> lines =
sc.parallelize(Arrays.asList("pandas", "i like
pandas"));
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

49

## Two ways to create RDDs

### 2. Loading data from external storage

- Python

```
lines = sc.textFile("/path/to/README.md")
```

- Scala

```
val lines = sc.textFile("/path/to/README.md")
```

- Java

```
JavaRDD<String> lines =
sc.textFile("/path/to/README.md");
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

50

# Resilient Distributed Dataset – RDD

- RDD Basics
- Creating RDDs
- **RDD Operations**
- Common Transformation and Actions
- Persistence (Caching)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

51

## RDD Operations

- Two types of operations
  - Transformations: operations that **return a new RDDs** e.g., `map()`, `filter()`
  - Actions: operations that return a **result** to the driver program or write it to storage such as `count()`, `first()`
- Treated differently by Spark
  - Transformation: lazy evaluation
  - Action: execution at any time



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

52

# Transformation

- Example 1. Use `filter()`

- Python

```
inputRDD = sc.textFile("log.txt")
errorsRDD = inputRDD.filter(lambda x: "error" in x)
```

- Scala

```
val inputRDD = sc.textFile("log.txt")
val errorsRDD = inputRDD.filter(line =>
    line.contains("error"))
```

- Java

```
JavaRDD<String> inputRDD = sc.textFile("log.txt");
JavaRDD<String> errorsRDD = inputRDD.filter(
    new Function<String, Boolean>() {
        public Boolean call(String x) {
            return x.contains("error");
        }
    });

```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

53

# Transformation

- `filter()`

- does not change the existing `inputRDD`
- returns a pointer to an entirely new RDD
- `inputRDD` still can be reused

- `union()`

```
errorsRDD = inputRDD.filter(lambda x: "error" in x)
warningsRDD = inputRDD.filter(lambda x: "warning" in x)
badLinesRDD = errorsRDD.union(warningsRDD)
```

- transformations can operate on any number of input RDDs

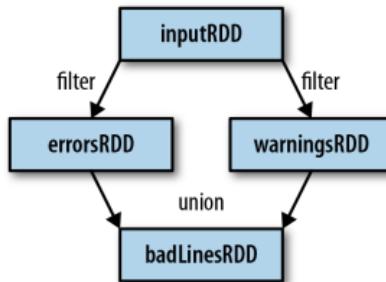


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

54

# Transformation

- Spark keeps track dependencies between RDDs, called the **lineage graph**
- Allow recovering lost data



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

55

# Actions

- Example. count the number of errors
- Python

```

print "Input had " + badLinesRDD.count() + " concerning lines"
print "Here are 10 examples:"
for line in badLinesRDD.take(10):
    print line
  
```

- Scala

```

println("Input had " + badLinesRDD.count() + " concerning lines")
println("Here are 10 examples:")
badLinesRDD.take(10).foreach(println)
  
```

- Java

```

System.out.println("Input had " + badLinesRDD.count() + " concerning
lines")
System.out.println("Here are 10 examples:")
for (String line: badLinesRDD.take(10)) {
    System.out.println(line);
}
  
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

56

# Resilient Distributed Dataset – RDD

- RDD Basics
- Creating RDDs
- RDD Operations
- Common Transformation and Actions
- Persistence (Caching)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

57

## RDD Basics

Transformations	Actions
map flatMap filter sample union groupByKey reduceByKey join cache	reduce collect count save lookupKey ...



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

58

# Transformations

<i>transformation</i>	<i>description</i>
<b>map(func)</b>	return a new distributed dataset formed by passing each element of the source through a function <i>func</i>
<b>filter(func)</b>	return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true
<b>flatMap(func)</b>	similar to map, but each input item can be mapped to 0 or more output items (so <i>func</i> should return a Seq rather than a single item)
<b>sample(withReplacement, fraction, seed)</b>	sample a fraction <i>fraction</i> of the data, with or without replacement, using a given random number generator <i>seed</i>
<b>union(otherDataset)</b>	return a new dataset that contains the union of the elements in the source dataset and the argument
<b>distinct([numTasks])</b>	return a new dataset that contains the distinct elements of the source dataset



59

# Transformations

<i>transformation</i>	<i>description</i>
<b>groupByKey([numTasks])</b>	when called on a dataset of (K, V) pairs, returns a dataset of (K, Seq[V]) pairs
<b>reduceByKey(func, [numTasks])</b>	when called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function
<b>sortByKey([ascending], [numTasks])</b>	when called on a dataset of (K, V) pairs where K implements Ordered, returns a dataset of (K, V) pairs sorted by keys in ascending or descending order, as specified in the boolean <i>ascending</i> argument
<b>join(otherDataset, [numTasks])</b>	when called on datasets of type (K, V) and (K, W), returns a dataset of (K, (V, W)) pairs with all pairs of elements for each key
<b>cogroup(otherDataset, [numTasks])</b>	when called on datasets of type (K, V) and (K, W), returns a dataset of (K, Seq[V], Seq[W]) tuples – also called groupWith
<b>cartesian(otherDataset)</b>	when called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements)



VIEN CONG NGHE THONG TIN VA TRUYEN THONG

60

# Actions

<i>action</i>	<i>description</i>
<b>reduce(<i>func</i>)</b>	aggregate the elements of the dataset using a function <i>func</i> (which takes two arguments and returns one), and should also be commutative and associative so that it can be computed correctly in parallel
<b>collect()</b>	return all the elements of the dataset as an array at the driver program – usually useful after a filter or other operation that returns a sufficiently small subset of the data
<b>count()</b>	return the number of elements in the dataset
<b>first()</b>	return the first element of the dataset – similar to <i>take(1)</i>
<b>take(<i>n</i>)</b>	return an array with the first <i>n</i> elements of the dataset – currently not executed in parallel, instead the driver program computes all the elements
<b>takeSample(<i>withReplacement</i>, <i>fraction</i>, <i>seed</i>)</b>	return an array with a random sample of <i>num</i> elements of the dataset, with or without replacement, using the given random number generator seed



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

61

# Actions

<i>action</i>	<i>description</i>
<b>saveAsTextFile(<i>path</i>)</b>	write the elements of the dataset as a text file (or set of text files) in a given directory in the local filesystem, HDFS or any other Hadoop-supported file system. Spark will call <i>toString</i> on each element to convert it to a line of text in the file
<b>saveAsSequenceFile(<i>path</i>)</b>	write the elements of the dataset as a Hadoop SequenceFile in a given path in the local filesystem, HDFS or any other Hadoop-supported file system. Only available on RDDs of key-value pairs that either implement Hadoop's Writable interface or are implicitly convertible to Writable (Spark includes conversions for basic types like Int, Double, String, etc).
<b>countByKey()</b>	only available on RDDs of type (K, V). Returns a `Map` of (K, Int) pairs with the count of each key
<b>foreach(<i>func</i>)</b>	run a function <i>func</i> on each element of the dataset – usually done for side effects such as updating an accumulator variable or interacting with external storage systems

62

# Resilient Distributed Dataset – RDD

- RDD Basics
- Creating RDDs
- RDD Operations
- Common Transformation and Actions
- Persistence (Caching)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

63

## Persistence levels

Level	Space used	CPU time	In memory	On disk	Comments
MEMORY_ONLY	High	Low	Y	N	
MEMORY_ONLY_SER	Low	High	Y	N	
MEMORY_AND_DISK	High	Medium	Some	Some	Spills to disk if there is too much data to fit in memory.
MEMORY_AND_DISK_SER	Low	High	Some	Some	Spills to disk if there is too much data to fit in memory. Stores serialized representation in memory.
DISK_ONLY	Low	High	N	Y	



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

64

# Persistence

- Example

```
val result = input.map(x => x * x)
result.persist(StorageLevel.DISK_ONLY)
println(result.count())
println(result.collect().mkString(", "))
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

65

## Books:

- Holden Karau, Andy Konwinski, Patrick Wendell & Matei Zaharia. Learning Spark. O'reilly
- TutorialsPoint. Spark Core Programming

Acknowledgement and References

## Slides:

- Paco Nathan. Intro to Apache Spark
- Harold Liu. Berkely Data Analytics Stack
- DataBricks. Intro to Spark Development



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

66

---

# Q&A



67



68



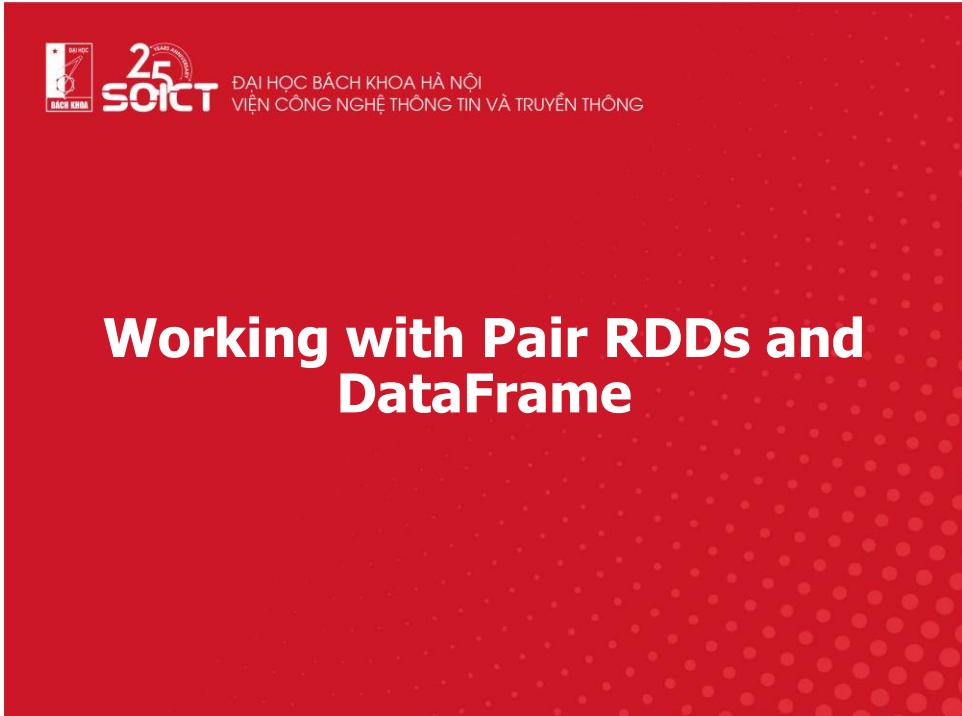
69



70



1



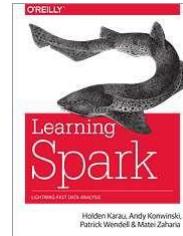
2

---

## From where to learn Spark ?



<http://spark.apache.org/>



<http://shop.oreilly.com/product/0636920028512.do>

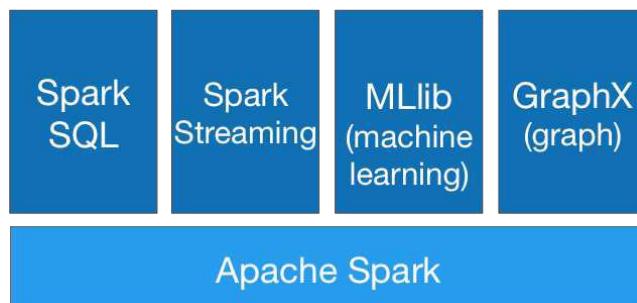


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

3

---

## Spark architecture



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

4

## Easy ways to run Spark ?

- ★ your IDE (ex. Eclipse or IDEA)
- ★ [Standalone Deploy Mode](#): simplest way to deploy Spark on a single machine
- ★ [Docker & Zeppelin](#)
- ★ [EMR](#)
- ★ Hadoop vendors ([Cloudera](#), [Hortonworks](#))



5

## Supported languages



6

# RDD (Resilient Distributed Dataset)

## RDD (Resilient Distributed Dataset)

- Resilient: If data in memory is lost, it can be recreated
  - Distributed: Processed across the cluster
  - Dataset: Initial data can come from a source such as a file, or it can be created programmatically
- RDDs are the fundamental unit of data in Spark**
- Most Spark programming consists of performing operations on RDDs**



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

8

## Creating RDD (I)

- Python
- `lines = sc.parallelize(["workshop", "spark"])`
- Scala
- `val lines = sc.parallelize(List("workshop", "spark"))`
- Java
- `JavaRDD<String> lines = sc.parallelize(NSArray.asList("workshop", "spark"))`



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

9

## Creating RDD (II)

Python

```
lines = sc.textFile("/path/to/file.txt")
```

Scala

```
val lines = sc.textFile("/path/to/file.txt")
```

Java

```
JavaRDD<String> lines = sc.textFile("/path/to/file.txt")
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

10

## RDD persistence

```
MEMORY_ONLY  
MEMORY_AND_DISK  
MEMORY_ONLY_SER  
MEMORY_AND_DISK_SER  
DISK_ONLY  
MEMORY_ONLY_2  
MEMORY_AND_DISK_2  
OFF_HEAP
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

11

## Working with RDDs



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

12

## RDDs

### **RDDs can hold any serializable type of element**

- Primitive types such as integers, characters, and booleans
- Sequence types such as strings, lists, arrays, tuples, and dicts (including nested data types)
- Scala/Java Objects (if serializable)
- Mixed types

### **§ Some RDDs are specialized and have additional functionality**

- Pair RDDs
- RDDs consisting of key-value pairs
- Double RDDs
- RDDs consisting of numeric data



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

13

# Creating RDDs from Collections

You can create RDDs from collections instead of files  
 –`sc.parallelize(collection)`

```
myData = ["Alice","Carlos","Frank","Barbara"]
> myRdd = sc.parallelize(myData)
> myRdd.take(2) ['Alice', 'Carlos']
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

14

# Creating RDDs from Text Files (1)

**For file-based RDDs, use `SparkContext.textFile`**

- Accepts a single file, a directory of files, a wildcard list of files, or a comma-separated list of files. Examples:

- `-sc.textFile("myfile.txt")`
- `-sc.textFile("mydata/")`
- `-sc.textFile("mydata/*.log")`
- `-sc.textFile("myfile1.txt,myfile2.txt")`

- Each line in each file is a separate record in the RDD

**Files are referenced by absolute or relative URI**

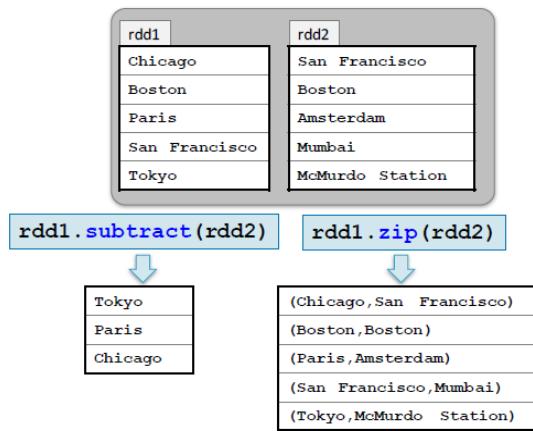
- Absolute URI:
- `-file:/home/training/myfile.txt`
- `-hdfs://nnhost/loudacre/myfile.txt`



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

15

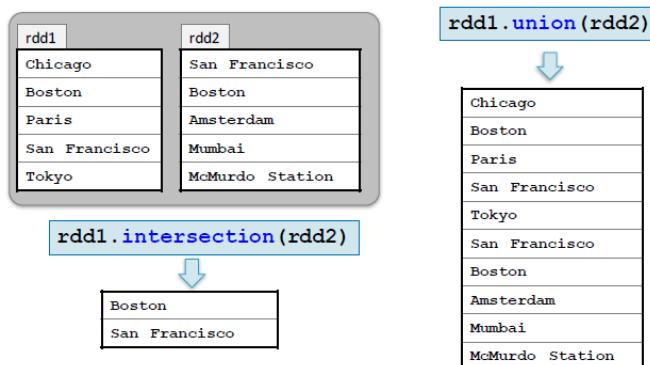
## Examples: Multi-RDD Transformations (1)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

16

## Examples: Multi-RDD Transformations (2)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

17

---

## Some Other General RDD Operations

### Other RDD operations

- first returns the first element of the RDD
- foreach applies a function to each element in an RDD
- top( $n$ ) returns the largest  $n$  elements using natural ordering

### Sampling operations

- sample creates a new RDD with a sampling of elements
- takeSample returns an array of sampled elements



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

18

---

## Other data structures in Spark

- ★ Paired RDD
- ★ DataFrame
- ★ DataSet

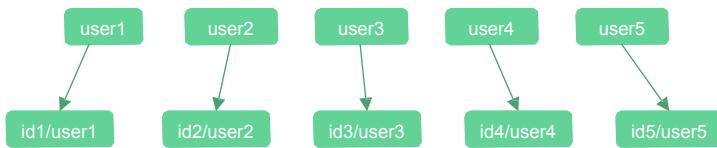


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

19

## Paired RDD

Paired RDD = an RDD of key/value pairs



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

20

Pair RDDs



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

21

## Pair RDDs

### § Pair RDDs are a special form of RDD

- Each element must be a key-value pair (a two-element *tuple*)
- Keys and values can be any type

Pair RDD

(key1, value1)
(key2, value2)
(key3, value3)
...

### § Why?

- Use with map-reduce algorithms
- Many additional functions are available for common data processing needs
- Such as sorting, joining, grouping, and counting



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

22

## Creating Pair RDDs

**The first step in most workflows is to get the data into key/value form**

- What should the RDD should be keyed on?
- What is the value?

**Commonly used functions to create pair RDDs**

- map
- flatMap / flatMapValues
- keyBy



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

23

## Example: A Simple Pair RDD

### Example: Create a pair RDD from a tab-separated file

```
> val users = sc.textFile(file).
  map(line => line.split('\t').
  map(fields => (fields(0), fields(1))))
```

user001\tFred Flintstone  
user090\tBugs Bunny  
user111\tHarry Potter  
...

(user001, Fred Flintstone)
(user090, Bugs Bunny)
(user111, Harry Potter)
...



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

24

## Example: Keying Web Logs by User ID

```
> sc.textFile(logfile).
  keyBy(line => line.split(' ') (2))
```

User ID  
 56.38.234.188 - 99788 "GET /KBDOC-00157.html HTTP/1.0" ...  
 56.38.234.188 - 99788 "GET /theme.css HTTP/1.0" ...  
 203.146.17.59 - 25254 "GET /KBDOC-00230.html HTTP/1.0" ...  
 ...

(99788, 56.38.234.188 - 99788 "GET /KBDOC-00157.html...")
(99788, 56.38.234.188 - 99788 "GET /theme.css...")
(25254, 203.146.17.59 - 25254 "GET /KBDOC-00230.html...")
...



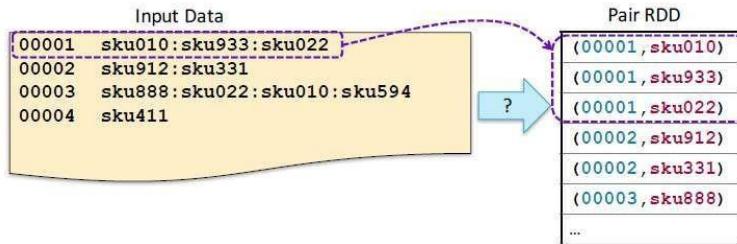
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

25

## Mapping Single Rows to Multiple Pairs

▪ How would you do this?

- Input: order numbers with a list of SKUs in the order
- Output: `order` (key) and `sku` (value)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

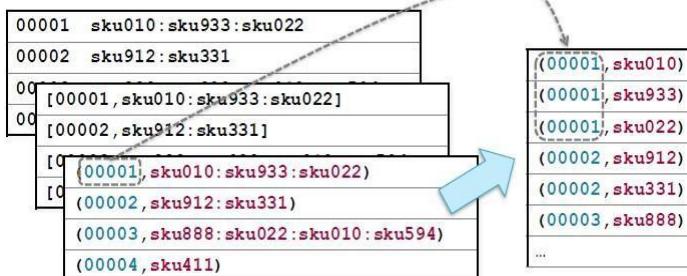
26

## Answer : Mapping Single Rows to Multiple Pairs

```

> sc.textFile(file) \
    .map(lambda line: line.split('\t')) \
    .map(lambda fields: (fields[0], fields[1])) \
    .flatMapValues(lambda skus: skus.split(':'))

```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

27

# Map-Reduce

- § **Map-reduce is a common programming model**
  - Easily applicable to distributed processing of large data sets
  
- § **Hadoop MapReduce is the major implementation**
  - Somewhat limited
  - Each job has one map phase, one reduce phase
  - Job output is saved to files
  
- § **Spark implements map-reduce with much greater flexibility**
  - Map and reduce functions can be interspersed
  - Results can be stored in memory
  - Operations can easily be chained



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

28

# Map-Reduce in Spark

- § **Map-reduce in Spark works on pair RDDs**

- § **Map phase**
  - Operates on one record at a time
  - “Maps” each record to zero or more new records
  - Examples: **map, flatMap, filter, keyBy**

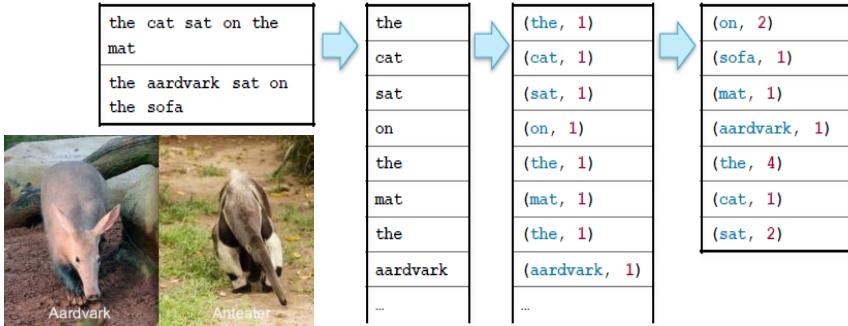
- § **Reduce phase**
  - Works on map output
  - Consolidates multiple records
  - Examples: **reduceByKey, sortByKey, mean**



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

29

## Example: Word Count



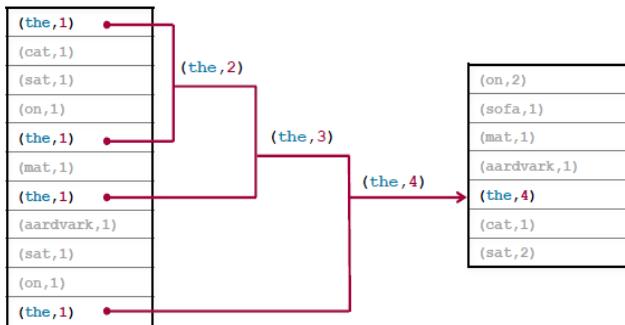
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

30

## reduceByKey

The function passed to `reduceByKey` combines values from two keys

- Function must be binary



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

31

```
> val counts = sc.textFile (file) . flatMap
  (line => line.split (' ')) . map (word => (word
,1)) . reduceByKey (v1 ,v2) => v1+v2)
```

OR

```
> val counts = sc.textFile (file) . flatMap
  (_.split (' ')) .
  map ((_,1)) .
  reduceByKey (_+_)
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

32

## Pair RDD Operations

§ In addition to map and reduceByKey operations, Spark has several operations specific to pair RDDs

### § Examples

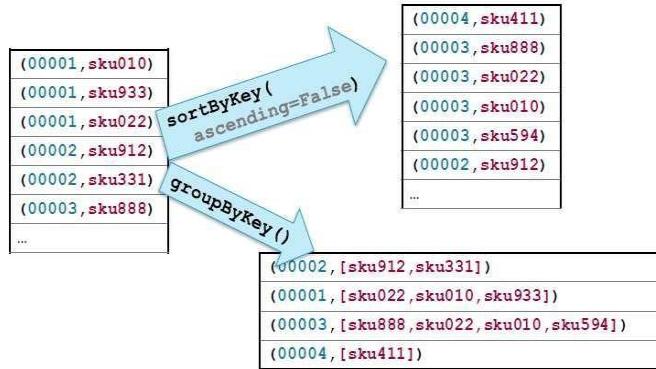
- **countByKey** returns a map with the count of occurrences of each key
- **groupByKey** groups all the values for each key in an RDD
- **sortByKey** sorts in ascending or descending order
- **join** returns an RDD containing all pairs with matching keys from two RDD



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

33

## Example: Pair RDD Operations

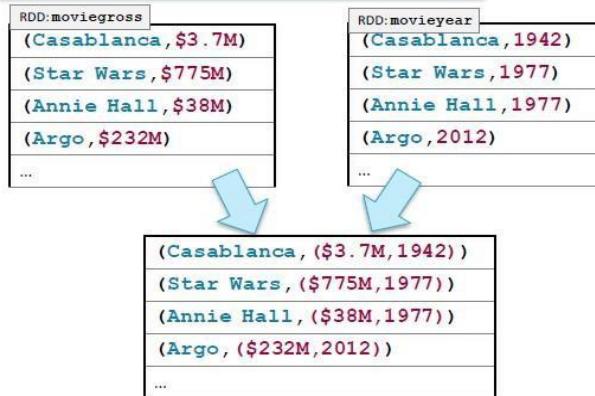


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

34

## Example: Joining by Key

```
> movies = moviegross.join(movieyear)
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

35

# Other Pair Operations

## § Some other pair operations

- keys** returns an RDD of just the keys, without the values
- values** returns an RDD of just the values, without keys
- lookup(key)** returns the value(s) for a key
- leftOuterJoin, rightOuterJoin , fullOuterJoin** join two RDDs, including keys defined in the left, right or either RDD respectively
- mapValues, flatMapValues** execute a function on just the values, keeping the key the same



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

36

## DataFrames and Apache Spark SQL



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

37

# What is Spark SQL?

## § What is Spark SQL?

- Spark module for structured data processing
- Replaces Shark (a prior Spark module, now deprecated)
- Built on top of core Spark

## § What does Spark SQL provide?

- The DataFrame API—a library for working with data as tables
- Defines DataFrames containing rows and columns
- DataFrames are the focus of this chapter!
- Catalyst Optimizer—an extensible optimization framework
- A SQL engine and command line interface



38

# SQL Context

## § The main Spark SQL entry point is a SQL context object

- Requires a **SparkContext** object
- The SQL context in Spark SQL is similar to Spark context in core Spark

## § There are two implementations

- SQLContext**
- Basic implementation
- HiveContext**
- Reads and writes Hive/HCatalog tables directly
- Supports full HiveQL language
- Requires the Spark application be linked with Hive libraries
- Cloudera recommends using **HiveContext**



39

# Creating a SQL Context

- § The Spark shell creates a `HiveContext` instance automatically
  - Call `sqlContext`
  - You will need to create one when writing a Spark application
  - Having multiple SQL context objects *is allowed*

- § A SQL context object is created based on the Spark context

Language: Scala

```
import org.apache.spark.sql.hive.HiveContext
val sqlContext = new HiveContext(sc)
import sqlContext.implicits._
```

40

# DataFrames

- § DataFrames are the main abstraction in Spark SQL

- Analogous to RDDs in core Spark
- A distributed collection of structured data organized into Named columns
- Built on a *base RDD* containing `Row` objects

## Creating a DataFrame from a Data Source

§ `sqlContext.read` returns a `DataFrameReader` object

§ `DataFrameReader` provides the functionality to load data into a `DataFrame`

§ Convenience functions

- `json(filename)`
- `parquet(filename)`
- `orc(filename)`
- `table(hive-tablename)`
- `jdb(url,table,options)`



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

42

## Example: Creating a DataFrame from a JSON File

```
Language: Scala
val sqlContext = new HiveContext(sc)
import sqlContext.implicits._
val peopleDF = sqlContext.read.json("people.json")
```

File: people.json

```
{"name": "Alice", "pcode": "94304"}
{"name": "Brayden", "age": 30, "pcode": "94304"}
 {"name": "Carla", "age": 19, "pcode": "10036"}
 {"name": "Diana", "age": 46}
 {"name": "Etienne", "pcode": "94104"}
```

age	name	pcode
null	Alice	94304
30	Brayden	94304
19	Carla	10036
46	Diana	null
null	Étienne	94104



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

43

## Example: Creating a DataFrame from a Hive/Impala Table

```
Language: Scala
val sqlContext = new HiveContext(sc)
import sqlContext.implicits._
val customerDF = sqlContext.read.table("customers")
```

Table: customers

cust_id	name	country
001	Ani	us
002	Bob	ca
003	Carlos	mx
...	...	...



cust_id	name	country
001	Ani	us
002	Bob	ca
003	Carlos	mx
...	...	...



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

44

## Loading from a Data Source Manually

### § You can specify settings for the DataFrameReader

- format**: Specify a data source type
- option**: A key/value setting for the underlying data source
- schema**: Specify a schema instead of inferring from the data source

### § Then call the generic base function load

```
sqlContext.read.
  format("com.databricks.spark.avro").
  load("/loudacre/accounts_avro")
```

```
sqlContext.read.
  format("jdbc").
  option("url","jdbc:mysql://localhost/loudacre").
  option("dbtable","accounts").
  option("user","training").
  option("password","training").
  load()
```



45

# Data Sources

## § Spark SQL 1.6 built-in data source types

- table
- json
- parquet
- jdbc
- orc

## § You can also use third party data source libraries, such as

- Avro (included in CDH)
- HBase
- CSV
- MySQL
- and more being added all the time



46

# DataFrame Basic Operations

- § Basic operations deal with DataFrame metadata (rather than its data)
- § Some examples
  - -schema returns a schema object describing the data
  - -printSchema displays the schema as a visual tree
  - -cache / persist persists the DataFrame to disk or memory
  - -columns returns an array containing the names of the columns
  - -dtypes returns an array of (column name,type) pairs
  - -explain prints debug information about the DataFrame to the console



47

# DataFrame Basic Operations

Language: Scala

```
> val peopleDF = sqlContext.read.json("people.json")
> peopleDF.dtypes.foreach(println)
(age, LongType)
(name, StringType)
(pcode, StringType)
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

48

# DataFrame Actions

- § Some DataFrame actions
  - –collect returns all rows as an array of Row objects
  - –take(*n*) returns the first *n* rows as an array of Row objects
  - –count returns the number of rows
  - –show(*n*) displays the first *n* rows  
(default=20)

Language: Scala

```
> peopleDF.count()
res7: Long = 5

> peopleDF.show(3)
age   name    pcode
null  Alice   94304
30   Brayden 94304
19   Carla   10036
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

49

# DataFrame Queries

## § DataFrame query methods return new DataFrames

- Queries can be chained like transformations

## § Some query methods

- distinct** returns a new DataFrame with distinct elements of this DF
- join** joins this DataFrame with a second DataFrame
- Variants for inside, outside, left, and right joins
- limit** returns a new DataFrame with the first **n** rows of this DF
- select** returns a new DataFrame with data from one or more columns of the base DataFrame
- where** returns a new DataFrame with rows meeting specified query criteria (alias for **filter**)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

50

# DataFrame Query Strings

## ■ Some query operations take strings containing simple query expressions

- Such as **select** and **where**

## ■ Example: **select**

age	name	pcode
null	Alice	94304
30	Brayden	94304
19	Carla	10036
46	Diana	null
null	Étienne	94104

peopleDF.  
select("age")

age
null
30
19
46
null

peopleDF.  
select("name", "age")

name	age
Alice	null
Brayden	30
Carla	19
Diana	46
Étienne	null



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

51

## Querying DataFrames using Columns

§ Columns can be referenced in multiple ways

The diagram illustrates how to select a single column from a DataFrame. It shows two Scala code snippets:

```
val ageDF = peopleDF.select(peopleDF("age"))
val ageDF = peopleDF.select($"age")
```

Both snippets result in a smaller DataFrame containing only the 'age' column:

age
null
30
19
46
null



52

## Joining DataFrames

§ A basic inner join when join column is in both DataFrames

The diagram shows two DataFrames being joined:

- peopleDF:**

age	name	pcode
null	Alice	94304
30	Brayden	94304
19	Carla	10036
46	Diana	null
null	Etienne	94104
- pcodesDF:**

pcodes	city	state
10036	New York	NY
87501	Santa Fe	NM
94304	Palo Alto	CA
94104	San Francisco	CA

The resulting joined DataFrame is:

pcodes	age	name	city	state
94304	null	Alice	Palo Alto	CA
94304	30	Brayden	Palo Alto	CA
10036	19	Carla	New York	NY
94104	null	Etienne	San Francisco	CA

Language: Python/Scala

```
peopleDF.join(pcodesDF, "pcodes")
```



53

# Joining DataFrames

- Specify type of join as inner (default), outer, left\_outer, right\_outer, or leftsemi

The diagram illustrates a join operation between two DataFrames. On the left, there is a DataFrame with columns: age, name, and pcode. The pcode column is highlighted with a red box. On the right, there are two code snippets: one in Python and one in Scala, both using the `join` method with the argument `"left_outer"`. Below the snippets, the resulting DataFrame is shown, which contains all columns from both original DataFrames. The pcode column from the first DataFrame is now aligned with the pcode column from the second DataFrame.

age	name	pcode
null	Alice	94304
30	Brayden	94304
19	Carla	10036
46	Diana	null
null	Étienne	94104

pcode	city	state
10036	New York	NY
87501	Santa Fe	NM
94304	Palo Alto	CA
94104	San Francisco	CA

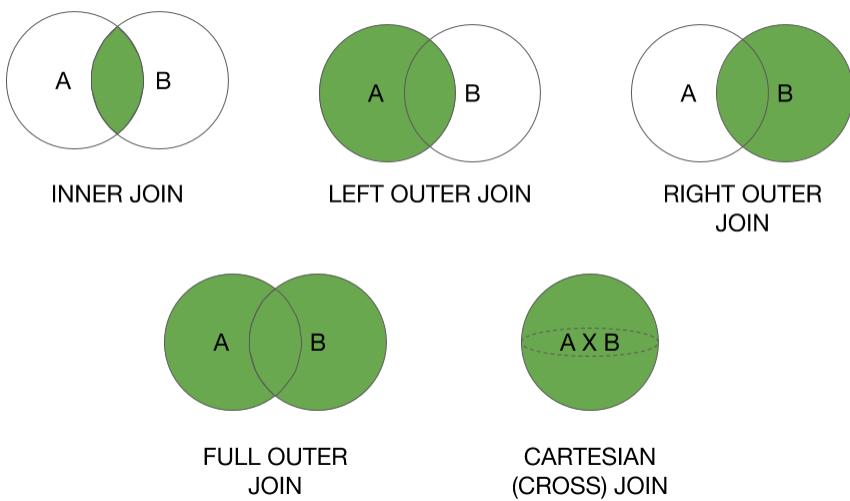
  

pcode	age	name	city	state
94304	null	Alice	Palo Alto	CA
94304	30	Brayden	Palo Alto	CA
10036	19	Carla	New York	NY
null	46	Diana	null	null
94104	null	Étienne	San Francisco	CA



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

54



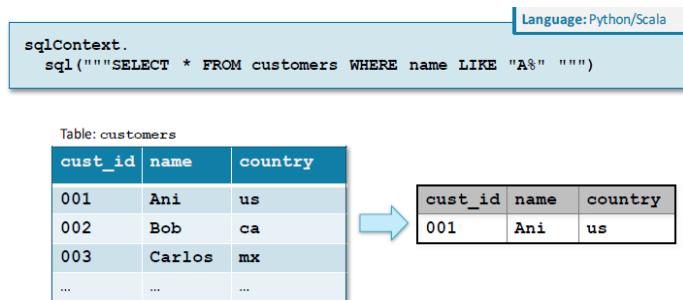
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

55

# SQL Queries

§ When using HiveContext, you can query Hive/Impala tables using HiveQL

- Returns a DataFrame



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

56

## Saving DataFrames

- § Data in DataFrames can be saved to a data source
- § Use DataFrame.write to create a DataFrameWriter
- § DataFrameWriter provides convenience functions to externally save the data represented by a DataFrame
  - **jdbc** inserts into a new or existing table in a database
  - **json** saves as a JSON file
  - **parquet** saves as a Parquet file
  - **orc** saves as an ORC file
  - **text** saves as a text file (string data in a single column only)
  - **saveAsTable** saves as a Hive/Impala table (**HiveContext** only)



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

57

## Options for Saving DataFrames

- § DataFrameWriter option methods
  - –format specifies a data source type
  - –mode determines the behavior if file or table already exists:
  - overwrite, append, ignore or error (default is error)
  - –partitionBy stores data in partitioned directories in the form
  - *column=value* (as with Hive/Impala partitioning)
  - –options specifies properties for the target data source
  - –save is the generic base function to write the data

Language: Python/Scala

```
peopleDF.write.  
  format("parquet").  
  mode("append").  
  partitionBy("age").  
  saveAsTable("people")
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

58

## DataFrames and RDDs

### § DataFrames are built on RDDs

- Base RDDs contain **Row** objects
- Use **rdd** to get the underlying RDD

peopleRDD = peopleDF.rdd

peopleDF		
age	name	pcode
null	Alice	94304
30	Brayden	94304
19	Carla	10036
46	Diana	null
null	Etienne	94104

peopleRDD		
Row	[null, Alice, 94304]	
Row	[30, Brayden, 94304]	
Row	[19, Carla, 10036]	
Row	[46, Diana, null]	
Row	[null, Etienne, 94104]	



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

59

# DataFrames and RDDs

§ Row RDDs have all the standard Spark actions and transformations

- Actions: `collect`, `take`, `count`, and so on
- Transformations: `map`, `flatMap`, `filter`, and so on

§ Row RDDs can be transformed into pair RDDs to use map-reduce methods

§ DataFrames also provide convenience methods (such as `map`, `flatMap`, and `foreach`) for converting to RDDs



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

60

# Working with Row Objects

- Use `Array`-like syntax to return values with type `Any`
- `row(n)` returns element in the  $n$ th column
- `row.fieldIndex("age")` returns index of the `age` column
- Use methods to get correctly typed values
- `row.getAs[Long]("age")`
- Use type-specific `get` methods to return typed values
- `row.getString(n)` returns  $n$ th column as a string
- `row.getInt(n)` returns  $n$ th column as an integer
- And so on



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

61

## Prerequisites

- Docker
- Zeppelin Docker Container
- Terminal Tools (Command Prompt, PowerShell)

Hanoi University of Science and Technology



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

64

64

## Working with Spark RDDs, Pair-RDDs



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

Hanoi University of Science and Technology

65

65

## RDD Operations

### Transformations

map()  
flatMap()  
filter()  
union()  
intersection()  
distinct()  
groupByKey()  
reduceByKey()  
sortByKey()  
join()

### Actions

count()  
collect()  
first(), top(n)  
take(n), takeOrdered(n)  
countByValue()  
reduce()  
foreach()  
...



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

66

66

## Lambda Expression

### PySpark WordCount example:

```
input_file = sc.textFile("/path/to/text/file")
map = input_file.flatMap(lambda line: line.split(" "))
.map(lambda word: (word, 1))
counts = map.reduceByKey(lambda a, b: a + b)
counts.saveAsTextFile("output")
```

**lambda arguments: expression**



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

67

67

## PySpark RDD API

<https://spark.apache.org/docs/latest/api/python/pyspark.html#pyspark.RDD>

`map(f, preservesPartitioning=False)`

[source]

Return a new RDD by applying a function to each element of this RDD.

```
>>> rdd = sc.parallelize(["b", "a", "c"])
>>> sorted(rdd.map(lambda x: (x, 1)).collect())
[('a', 1), ('b', 1), ('c', 1)]
```

`flatMap(f, preservesPartitioning=False)`

[source]

Return a new RDD by first applying a function to all elements of this RDD, and then flattening the results.

```
>>> rdd = sc.parallelize([2, 3, 4])
>>> sorted(rdd.flatMap(lambda x: range(1, x)).collect())
[1, 1, 1, 2, 2, 3]
>>> sorted(rdd.flatMap(lambda x: [(x, x), (x, x)]).collect())
[(2, 2), (2, 2), (3, 3), (3, 3), (4, 4), (4, 4)]
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

68

68

## Practice with flight data (1)

Data: **airports.dat** (<https://openflights.org/data.html>)

[*Airport ID, Name, City, Country, IATA, ICAO, Latitude, Longitude, Altitude, Timezone, DST, Tz database, Type, Source*]

**Try to do somethings:**

- Create RDD from textfile
- Count the number of airports
- Filter by country
- Group by country
- Count the number of airports in each country



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

69

69

## Practice with flight data (2)

- Data: **airports.dat** (<https://openflights.org/data.html>)  
 $[Airport\ ID, Name, City, Country, IATA, ICAO, Latitude, Longitude, Altitude, Timezone, DST, Tz\ database, Type, Source]$
- Data: **routes.dat**  
 $[Airline, Airline\ ID, Source\ airport, Source\ airport\ ID, Destination\ airport, Destination\ airport\ ID, Codeshare, Stops, Equipment]$

**Try to do somethings:**

- Join 2 RDD
- Count the number of flights arriving in each country



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

70

70

## Working with DataFrame and Spark SQL



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

Hanoi University of Science and Technology

71

71

## Creating a DataFrame(1)

```
%pyspark
from pyspark.sql import *

Employee = Row("firstName", "lastName", "email", "salary")

employee1 = Employee('Basher', 'armbrust', 'bash@edureka.co', 100000)
employee2 = Employee('Daniel', 'meng', 'daniel@stanford.edu', 120000 )
employee3 = Employee('Muriel', None, 'muriel@waterloo.edu', 140000 )
employee4 = Employee('Rachel', 'wendell', 'rach_3@edureka.co', 160000 )
employee5 = Employee('Zach', 'galifianakis', 'zach_g@edureka.co', 160000 )

employees = [employee1, employee2, employee3, employee4, employee5]

print(Employee[0])

print(employees)

dframe = spark.createDataFrame(employees)
dframe.show()
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

72

72

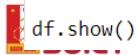
## Creating a DataFrame

### From CSV file:

```
%pyspark
flightData2015 = spark\
.read\
.option("inferSchema", "true")\
.option("header", "true")\
.csv("/usr/zeppelin/module9/2015-summary.csv")

flightData2015.show()
```

```
%pyspark
from pyspark.sql import *
list = [('Ankit', 25), ('Jalfaizy', 22), ('saurabh', 20), ('Bala', 26)]
rdd = sc.parallelize(list)
people = rdd.map(lambda x: Row(name=x[0], age=int(x[1])))
df = spark.createDataFrame(people)
```



73

73

## DataFrame APIs

- **DataFrame**: show(), collect(), createOrReplaceTempView(), distinct(), filter(), select(), count(), groupBy(), join()...
- **Column**: like()
- **Row**: row.key, row[key]
- **GroupedData**: count(), max(), min(), sum(), ...

<https://spark.apache.org/docs/latest/api/python/pyspark.sql.html>

Hanoi University of Science and Technology



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

74

74

## Spark SQL

- Create a temporary view
- Query using SQL syntax

```
%pyspark
flightData2015.createOrReplaceTempView("flight_data_2015")

maxSql = spark.sql("""
SELECT DEST_COUNTRY_NAME, sum(count) as destination_total
FROM flight_data_2015
GROUP BY DEST_COUNTRY_NAME
ORDER BY sum(count) DESC
LIMIT 5
""")

maxSql.show()
```



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

75

75



76



77



78

# Lecture 5: Data Preparation

## Data Preparation

- Introduction to Data Preparation.
- Types of Data.
- Outliers.
- Data Transformation.
- Missing Data.

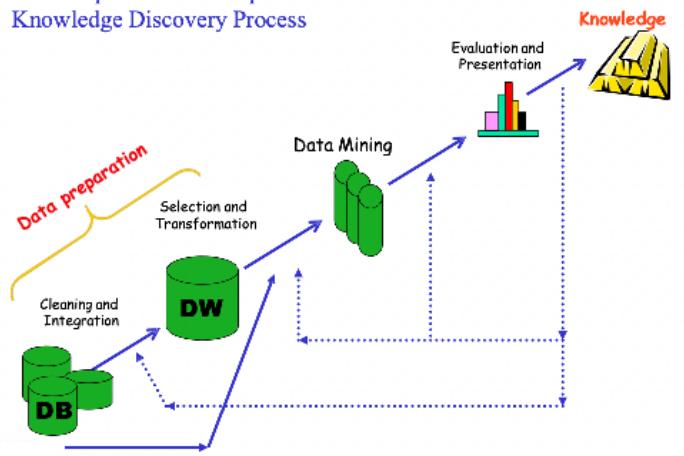
## Why Prepare Data

- Some data preparation is needed for all mining tools.
- The purpose of preparation is to transform datasets so that their information content is best exposed to the mining tool.
- Error prediction rate should be lower (or the same) after the preparation as before it.
- Preparing data also prepares the miner so that when using prepared data the miner produces better models, faster.
- Good data is a prerequisite for producing effective models of any type.
- Data need to be formatted for a given software tool.
- Data need to be made adequate for a given method.
- Data in the real world is dirty.
  - o **incomplete:** lacking attribute values, lacking certain attributes of interest, or containing only aggregate data.
    - e.g., occupation="".
  - o **noisy:** containing errors or outliers.
    - e.g., Salary = “-10”, Age = “222”.
  - o **inconsistent:** containing discrepancies in codes or names.
    - e.g., Age= “42”, Birthday = “03/07/1997”.
    - e.g., Was rating “1,2,3”, now rating “A, B, C”.
    - e.g., discrepancy between duplicate records.

## Major Tasks in Data Preparation

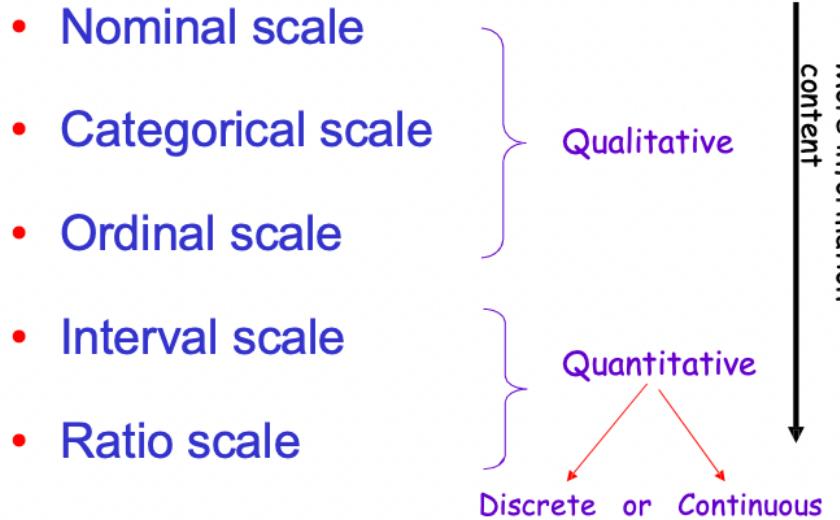
- Data discretization:
  - o Part of data reduction but with particular importance, especially for numerical data.
- Data cleaning:
  - o Fill in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies.
- Data integration:
  - o Integration of multiple databases, data cubes, or files.
- Data transformation:
  - o Normalization and aggregation.
- Data reduction:
  - o Obtains reduced representation in volume but produces the same or similar analytical results.

Data Preparation as a step in the Knowledge Discovery Process



## Types of Data

## Types of Measurements



## Types of Measurements: Example

- Nominal:
  - ID numbers, Names of people
- Categorical:
  - eye color, zip codes
- Ordinal:
  - rankings (e.g., taste of potato chips on a scale from 1-10), grades, height in {tall, medium, short}
- Interval:
  - calendar dates, temperatures in Celsius or Fahrenheit, GRE (Graduate Record Examination) and IQ scores
- Ratio:
  - temperature in Kelvin, length, time, counts

## Data Conversion

- Some tools can deal with nominal values but other need fields to be numeric.
- Convert ordinal fields to numeric to be able to use ">" and "<" comparisons on such fields.
  - A → 4.0
  - A- → 3.7
  - B+ → 3.3
  - B → 3.0
- Multi-valued, unordered attributes with small no. of values:
  - E.g., Color=Red, Orange, Yellow, ..., Violet

- For each value  $v$  create a binary “flag” variable  $C_v$ , which is 1 if  $\text{Color}=v$ , 0 otherwise.

## Conversion: Nominal, Many Values

- Examples:
  - US State Code (50 values).
  - Profession Code (7,000 values, but only few frequent).
- Ignore ID-like fields whose values are unique for each record.
- For other fields, group values “naturally”:
  - E.g., 50 US States => 3 or 5 regions.
  - Profession – select most frequent ones, group the rest.
- Create binary flag-fields for selected values.

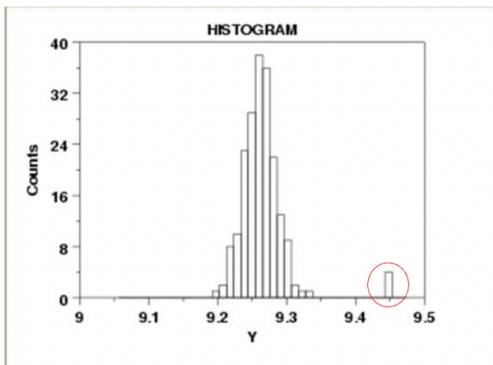
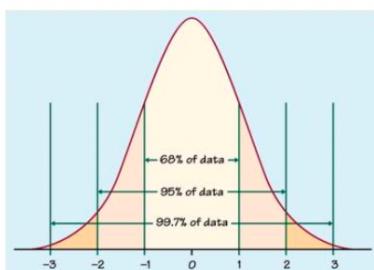
## Outliers

- Outliers are values thought to be out of range.
  - **“An outlier is an observation that deviates so much from other observations as to arouse suspicion that it was generated by a different mechanism”.**
  - Can be detected by standardizing observations and label the standardized values outside a predetermined bound as outliers.
  - Outlier detection can be used for fraud detection or data cleaning.
- Approaches:
  - Do nothing.
  - Enforce upper and lower bounds.
  - Let binning handle the problem.

## Outlier Detection

- **Univariate:**
  - **Compute mean and std. deviation.** For  $k=2$ ,  $k=3$ ,  $x$  is an outlier if outside limits (normal distribution assumed).

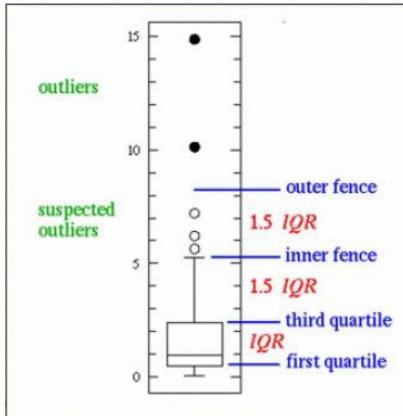
$$(x - ks, x + ks)$$



- **Boxplot:** An observation is an extreme outlier if:

$(Q1 - 3 \times IQR, Q3 + 3 \times IQR)$ , where  $IQR = Q3 - Q1$

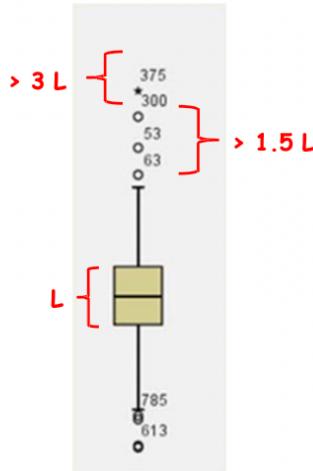
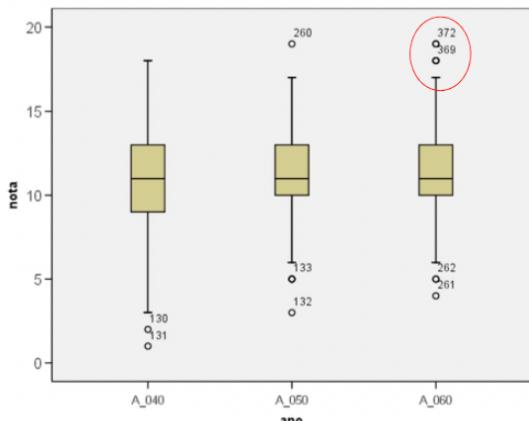
(*IQR = Inter Quartile Range*)



and declared a mild outlier if it lies outside of the interval

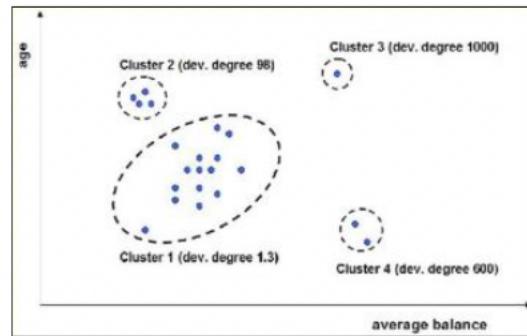
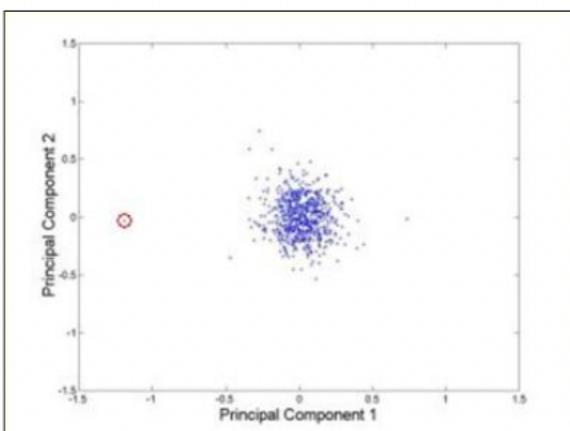
$(Q1 - 1.5 \times IQR, Q3 + 1.5 \times IQR)$ .

<http://www.physics.csbsju.edu/stats/box2.html>



- Multivariate:

- o Clustering:
  - Very small clusters are outliers.

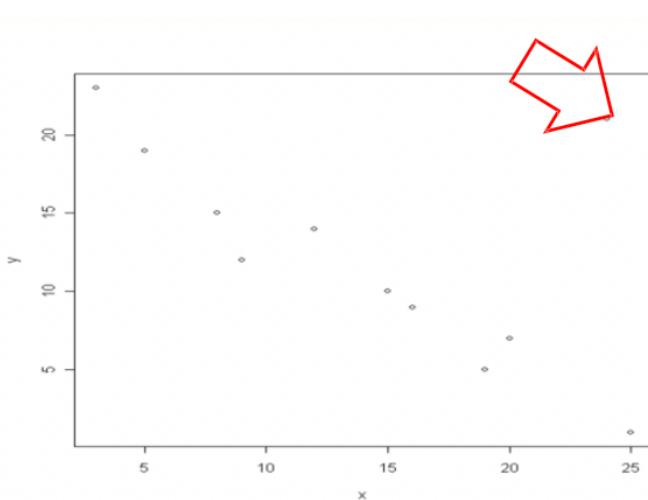
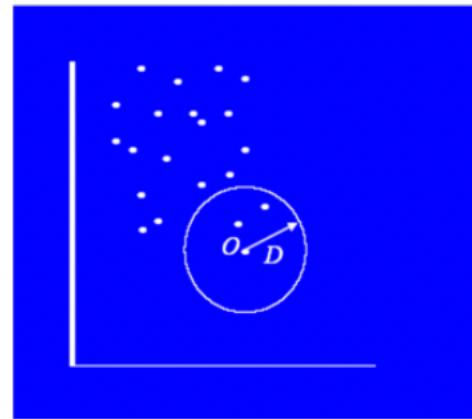
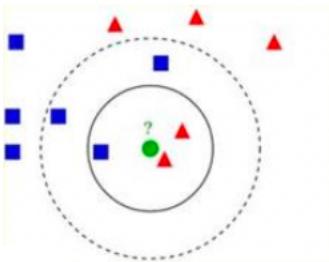


<http://www.ibm.com/developerworks/data/library/techarticle/dm-0811wurst/>

- o Distance based:

- An instance with very few neighbors within D is regarded as an outlier.

## Knn algorithm



A bi-dimensional outlier that is not an outlier in either of its projections.

## Data Transformation

### Normalization

- For distance-based methods, normalization helps to prevent that attribute with large ranges outweigh attributes with small ranges.
  - **min-max normalization**

$$v' = \frac{v - \min_v}{\max_v - \min_v} (\text{new\_max}_v - \text{new\_min}_v) + \text{new\_min}_v$$

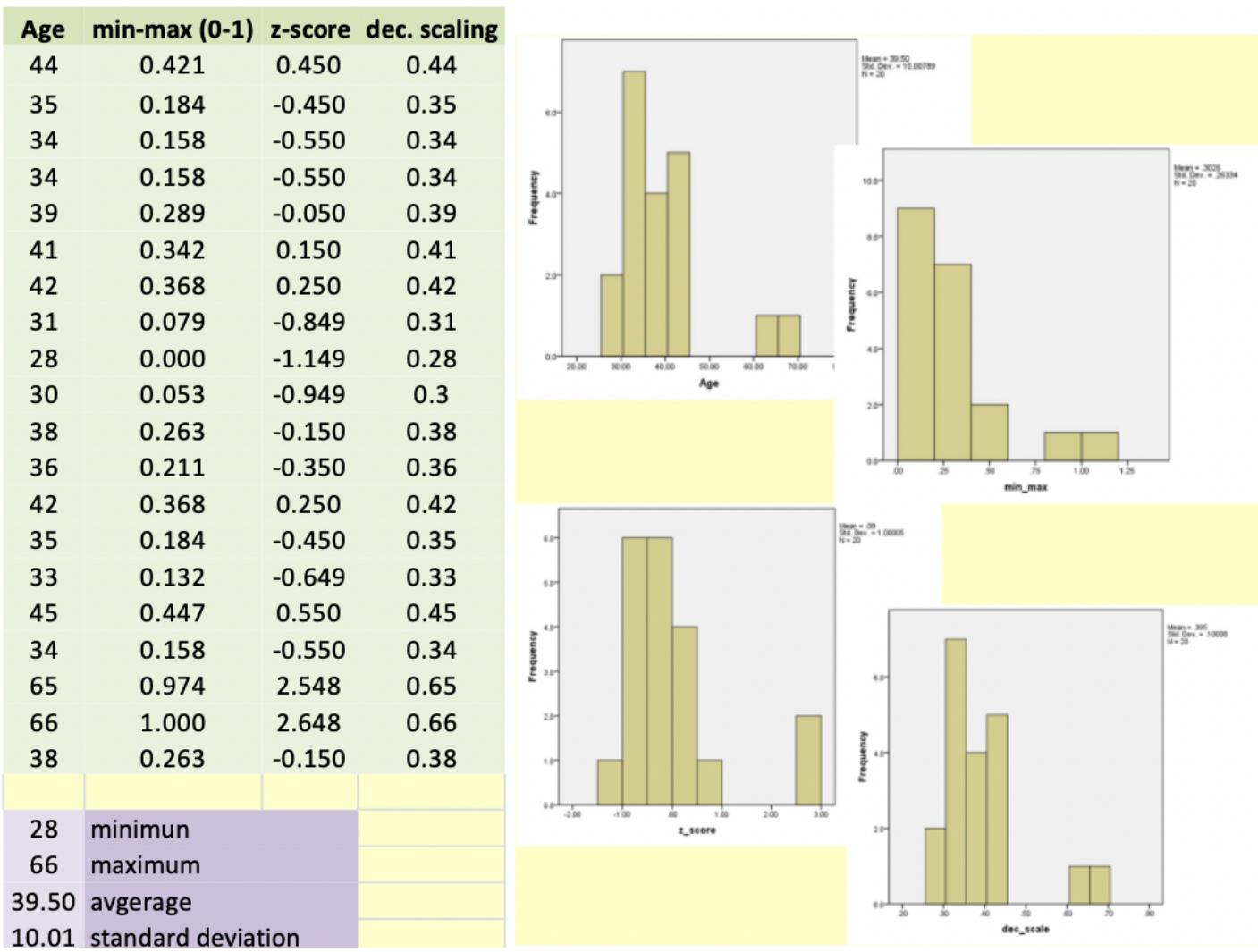
- **z-score normalization**

$$v' = \frac{v - \bar{v}}{\sigma_v} \quad \text{does not eliminate outliers}$$

- **normalization by decimal scaling**

$$v' = \frac{v}{10^j} \quad \text{Where } j \text{ is the smallest integer such that } \text{Max}(|v'|) < 1$$

range: -986 to 917 => j=3   -986 -> -0.986   917 -> 0.917



## Data Transformation

- It is the process to create new attributes:
  - o Often called transforming the attributes or the attribute set.
- Data transformation usually combines the original raw attributes using different mathematical formulas originated in business models or pure mathematical formulas.
- **Linear Transformations:**
  - o Normalizations may not be enough to adapt the data to improve the generated model.
  - o Aggregating the information contained in various attributes might be beneficial.
  - o If B is an attribute subset of the complete set A, A new attribute Z can be obtained by a linear combination:

$$Z = r_1 B_1 + r_2 B_2 + \cdots + r_m B_M$$

- **Quadratic Transformations:**
  - o In Quadratic Transformations, a new attribute is built as follows:

$$Z = r_{1,1} B_1^2 + r_{1,2} B_1 B_2 + \cdots + r_{m-1,m} B_{m-1} B_m + r_{m,m} B_m^2,$$

- These kinds of transformations have been thoroughly studied and can help to transform data to make it separable.
- **Non-polynomial Approximations of Transformations:**
  - Sometimes polynomial transformations are not enough.
  - For examples, guessing whether a set of triangles are congruent is not possible by simply observing their coordinates.
    - Computing the length of their segments will easily solve the problem → non polynomial transformation.

$$A = \sqrt{(X_1 - X_2)^2 + (Y_1 - Y_2)^2}$$

- **Polynomial Approximations of Transformations:**
  - We have observed that specific transformations may be needed to extract knowledge.
    - But help from an expert is not always available.
  - When no knowledge is available, a transformation  $f$  can be approximated via a polynomial transformation using a brute search with one degree at a time.
    - Using the Weistrass approximation, there is a polynomial function  $f$  that takes the value  $Y_i$  for each instance  $X_i$ .

$$Y = f(X_1, X_2, \dots, X_n)$$

- There are many polynomials verifying  $Y=f(X)$  as we want.
- As the number of instances in the data set increases, the approximations will be better.
- We can use computer assistance to approximate the intrinsic information.
- When the intrinsic transformation is polynomial, we need to add the cartesian product of the attributes needed for the polynomial degree approximation.
- Sometimes the approximation obtained must be rounded to avoid the limitations of the computer digital precision.
- **Rank transformation:**
  - A change in an attribute distribution can result in a change of model performance.
  - The simplest transformation to accomplish this in numerical attributes is to replace the value of an attribute with its rank.
  - The attribute will be transformed into a new attribute containing integer values ranging from 1 to  $m$ , being  $m$  the number of instances in the data set.
  - Next, we can transform the ranks to normal scores representing their probabilities in the normal distribution by spreading these values on the Gaussian curve using a simple transformation given by:

$$y = \Phi^{-1} \left( \frac{r_i - \frac{3}{8}}{m + \frac{1}{4}} \right)$$

- Being  $r_i$  the rank of the observation  $i$  and  $\Phi$  the cumulative normal function.
- Note: This transformation cannot be applied separately to the training and test partitions.
- **Box-Cox Transformations:**
  - When selecting the optimal transformation for an attribute is that we do not know in advance which transformation will be the best.
  - The Box-Cox transformation aims to transform a continuous variable into an almost normal distribution.
  - This can be achieved by mapping the values using following set of transformations:

$$y = \begin{cases} x^{\lambda-1}/\lambda, & \lambda \neq 0 \\ \log(x), & \lambda = 0 \end{cases}$$

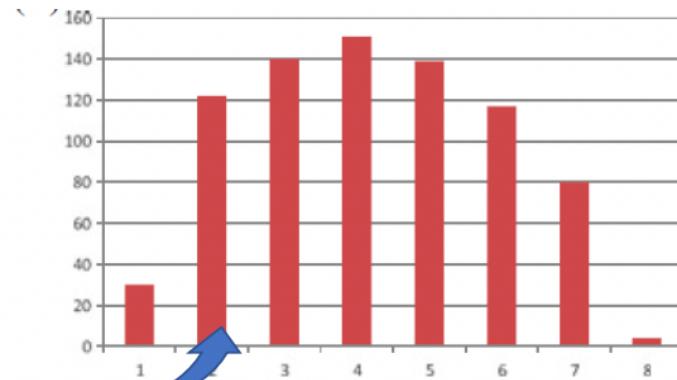
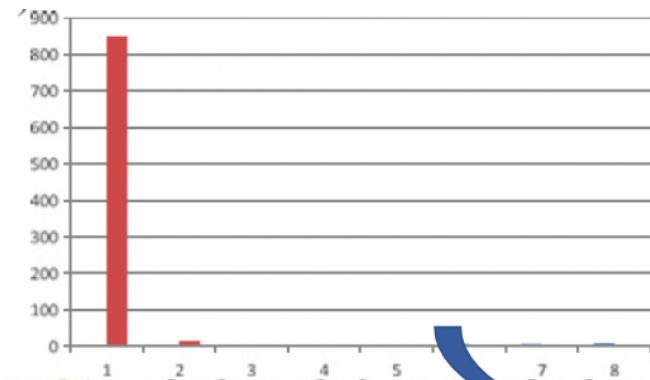
- All linear, inverse, quadratic and similar transformations are special cases of the Box-Cox transformations.
- Please note that all the values of variable  $x$  in the previous slide must be positive. If we have negative values in the attribute, we must add a parameter  $c$  to offset such negative values:

$$y = \begin{cases} (x + c)^{\lambda-1}/g\lambda, & \lambda \neq 0 \\ \log(x + c)/g, & \lambda = 0 \end{cases}$$

- The parameter  $g$  is used to scale the resulting values, and it is often considered as the geometric mean of the data.
- The value of  $\lambda$  is iteratively found by testing different values in the range from  $-3.0$  to  $3.0$  in small steps until the resulting attribute is as close as possible to the normal distribution.

#### - Spreading the Histogram:

- Spreading the histogram is a special case of Box-Cox transformations.
- As Box-Cox transforms the data to resemble a normal distribution, the histogram is thus spread as shown here.



- When the user is not interested in converting the distribution to a normal one, but just spreading it, we can use two special cases of Box-Cox transformations.
  - 1. Using the logarithm (with an offset if necessary) can be used to spread the right side of the histogram:  $y = \log(x)$ .
  - 2. If we are interested in spreading the left side of the histogram, we can simply use the power transformation  $y = x^{\lambda}$ .

#### - Nominal to Binary Transformation:

- The presence of nominal attributes in the data set can be problematic, especially if the Data Mining (DM) algorithm used cannot correctly handle them.
- The first option is to transform the nominal variable to a numeric one.
- Although simple, this approach has two big drawbacks that discourage it:
  - With this transformation we assume an ordering of the attribute values.
  - The integer values can be used in operations as numbers, whereas the nominal values cannot.
- To avoid the problems, a very typical transformation used for DM methods is to map each nominal attribute to a set of newly generated attributes.
- If  $N$  is the number of different values the nominal attribute has, we will substitute the nominal variable with a new set of binary attributes, each one representing one of the  $N$  possible values.

- For each instance, only one of the N newly created attributes will have a value of 1, while the rest will have the value of 0.
- This transformation is also referred in the literature as 1-to-N transformation.
- A problem with this kind of transformation appears when the original nominal attribute has a large cardinality:
  - The number of attributes generated will be large as well.
  - Resulting in a very sparse data set which will lead to numerical and performance problems.
- **Transformations via Data Reduction:**
  - When the data set is very large, performing complex analysis and DM can take a long computing time.
  - Data reduction techniques are applied in these domains to reduce the size of the data set while trying to maintain the integrity and the information of the original data set as much as possible.
  - Mining on the reduced data set will be much more efficient and it will also resemble the results that would have been obtained using the original data set.
  - The main strategies to perform data reduction are Dimensionality Reduction (DR) techniques.
  - They aim to reduce the number of attributes or instances available in the data set.
  - Well known attribute reduction techniques are Wavelet transforms or **Principal Component Analysis (PCA)**.
  - Many techniques can be found for reducing the dimensionality in the number of instances, like the use of clustering techniques, parametric methods and so on.
  - The use of binning and discretization techniques is also useful to reduce the dimensionality and complexity of the data set.
  - They convert numerical attributes into nominal ones, thus drastically reducing the cardinality of the attributes involved.

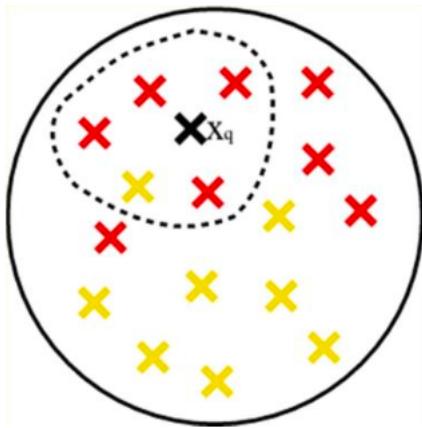
## Missing Data

- Data is not always available.
  - E.g., many tuples have no recorded value for several attributes, such as customer income in sales data.
- Missing data may be due to:
  - equipment malfunction.
  - inconsistent with other recorded data and thus deleted.
  - data not entered due to misunderstanding.
  - certain data may not be considered important at the time of entry.
  - not register history or changes of the data.
- Missing data may need to be inferred.
- Missing values may carry some information content.
  - E.g., a credit application may carry information by noting which field the applicant did not complete.
- There are always MVs in a real dataset.
  - MVs may have an impaction modelling, in fact, they can destroy it!
  - Some tools ignore missing values, others use some metric to fill in replacements.
  - The modeler should avoid default automated replacement techniques.
  - Difficult to know limitations, problems and introduced bias.
  - Replacing missing values without elsewhere capturing that information removes information from the dataset.

## How to Handle Missing Data?

- Ignore records (use only cases with all values)
  - o Usually done when class label is missing as most prediction methods do not handle missing data well.
  - o Not effective when the percentage of missing values per attribute varies considerably as it can lead to insufficient and/or biased sample sizes.
- Ignore attributes with missing values.
- Use only features (attributes) with all values (may leave out important features).
- Fill in the missing value manually.
  - o tedious + infeasible?
- Use a global constant to fill in the missing value.
  - o E.g., "unknown". (May create a new class).
- Use the attribute mean to fill in the missing value.
  - o It will do the least harm to the mean of existing data.
  - o If the mean is to be unbiased.
  - o What if the standard deviation is to be unbiased?
- Use the attribute mean for all samples belonging to the same class to fill in the missing value.
- Use the most probable value to fill in the missing value.
  - o Inference-based such as Bayesian formula or decision tree.
  - o Identify relationships among variables.
    - Linear regression, Multiple linear regression, Non-linear regression.
  - o Nearest-Neighbor estimator
    - Finding the k neighbors nearest to the point and fill in the most frequent value or the average value.
    - Finding neighbors in a large dataset may be slow.

## Nearest-neighbor



## How to Handle Missing Data?

- Note that, it is as important to avoid adding bias and distortion to the data as it is to make the information available.
  - o Bias is added when a wrong value is filled in.
- No matter what techniques you use to conquer the problem, it comes at a price. The more guessing you must do, the further away from the real data the database becomes. Thus, in turn, it can affect the accuracy and validation of the mining results.

## Summary

- Every real-world data set needs some kind of data pre-processing.
- Deal with missing values.
- Correct erroneous values.
- Select relevant attributes.
- Adapt data set format to the software tool to be used.
- In general, data pre-processing consumes more than 60% of a data mining project effort.

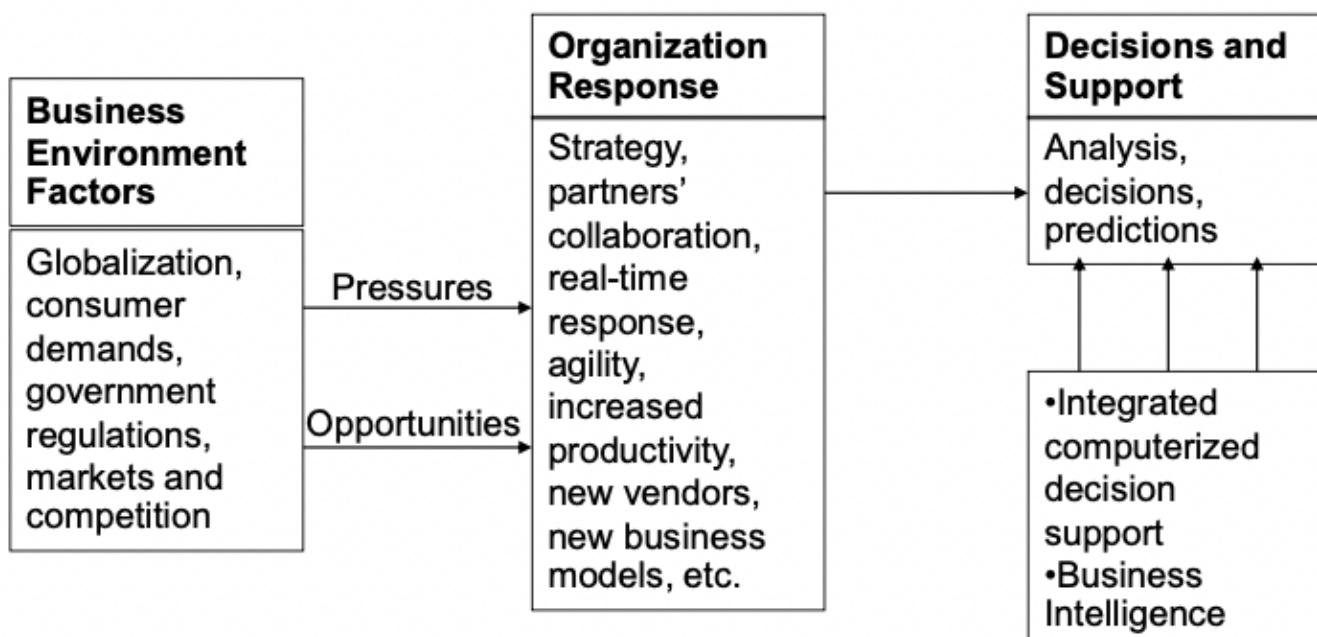
## Table of Contents

<b><i>Data Preparation</i></b> .....	<b>1</b>
<b><i>Why Prepare Data</i></b> .....	<b>1</b>
<b><i>Major Tasks in Data Preparation</i></b> .....	<b>1</b>
<b><i>Types of Data</i></b> .....	<b>2</b>
<b><i>Types of Measurements</i></b> .....	<b>2</b>
<b><i>Types of Measurements: Example</i></b> .....	<b>2</b>
<b><i>Data Conversion</i></b> .....	<b>2</b>
<b><i>Conversion: Nominal, Many Values</i></b> .....	<b>3</b>
<b><i>Outliers</i></b> .....	<b>3</b>
<b><i>Outlier Detection</i></b> .....	<b>3</b>
<b><i>Data Transformation</i></b> .....	<b>5</b>
<b><i>Normalization</i></b> .....	<b>5</b>
<b><i>Data Transformation</i></b> .....	<b>6</b>
<b><i>Missing Data</i></b> .....	<b>9</b>
<b><i>How to Handle Missing Data?</i></b> .....	<b>10</b>
<b><i>Nearest-neighbor</i></b> .....	<b>10</b>
<b><i>How to Handle Missing Data?</i></b> .....	<b>10</b>
<b><i>Summary</i></b> .....	<b>11</b>

# Lecture 6.1: Introduction to Decision Support Systems

## Decision Making

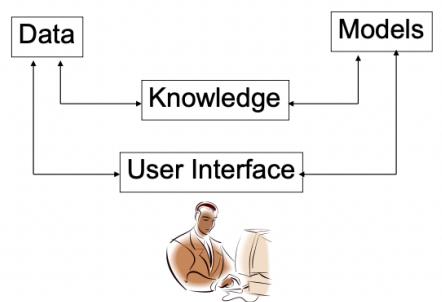
- Business Environment Factors.
  - o Markets: strong competition, global markets, market on Internet.
  - o Consumer demands: customization, quality, diversity, delivery.
  - o Technology: more innovations, more obsolescence rate, more information overload.
  - o Societal: more regulation and deregulation, more diversified workforce, more social responsibility.
- Business Pressure-Response-Support Model



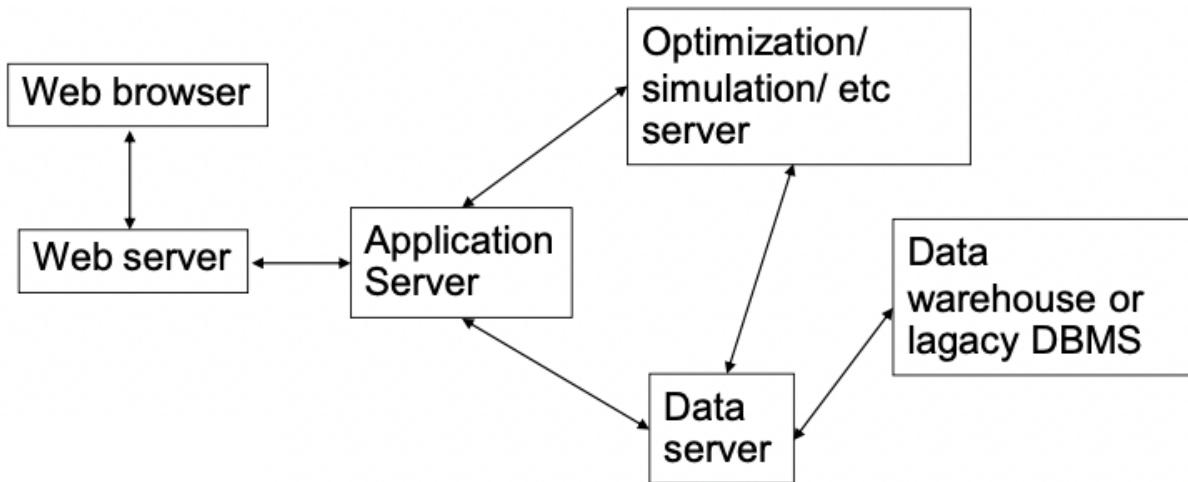
- Process of Decision Making:
  - o Define the problem (i.e., a decision situation that may deal with some difficulty or with an opportunity).
  - o Construct a model that describes the real-world problem.
  - o Identify possible solutions to the modeled problem and evaluate the solutions.

## Decision Support Systems

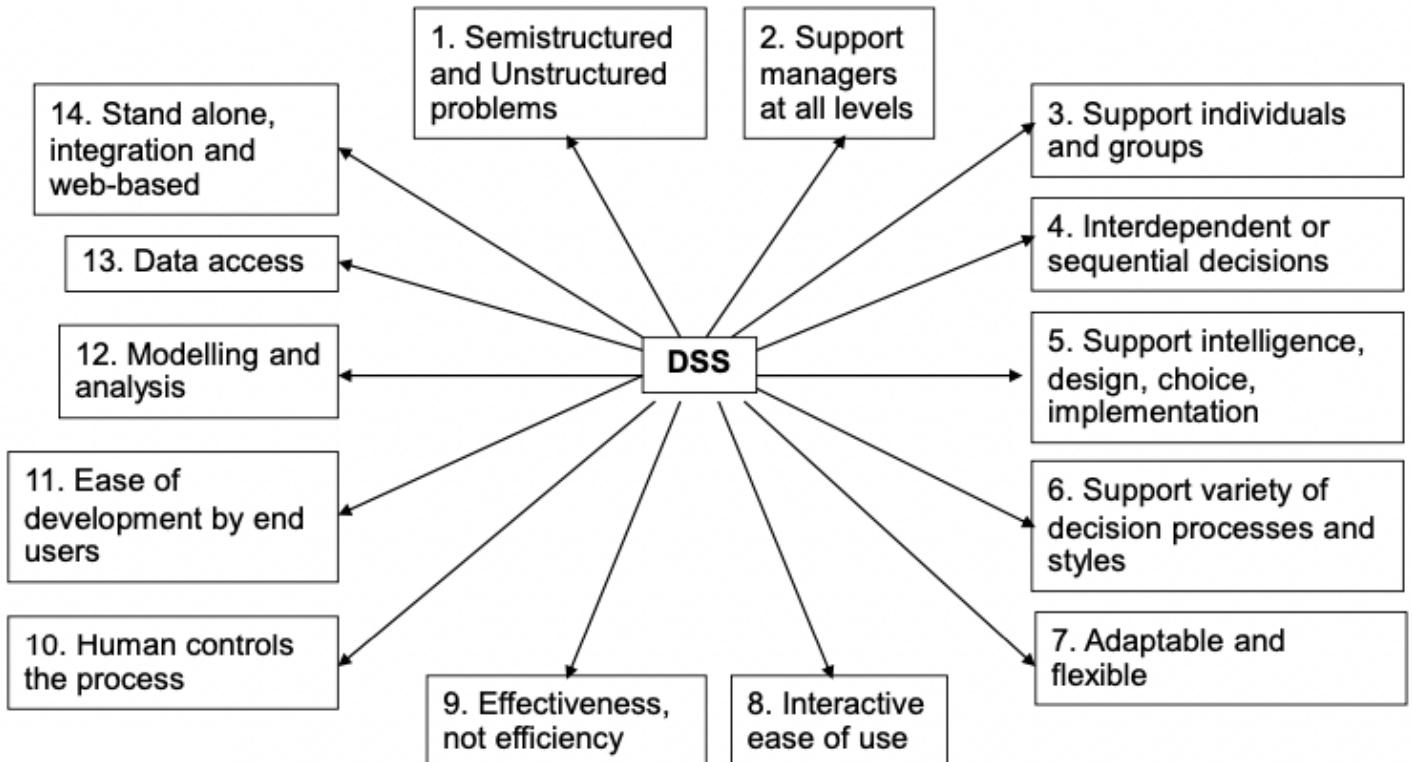
- Definition I (Keen and Scott-Morton):
  - o Decision Support Systems couple the intellectual resources of individuals with the capabilities of the computer to improve the quality of decisions.
  - o It is a computer-based support system for management decision makers who deal with semi-structured problems.
- High-level Architecture of DSS:



- Multitiered architecture for incorporating optimization, simulation, and other models into web based DSS:



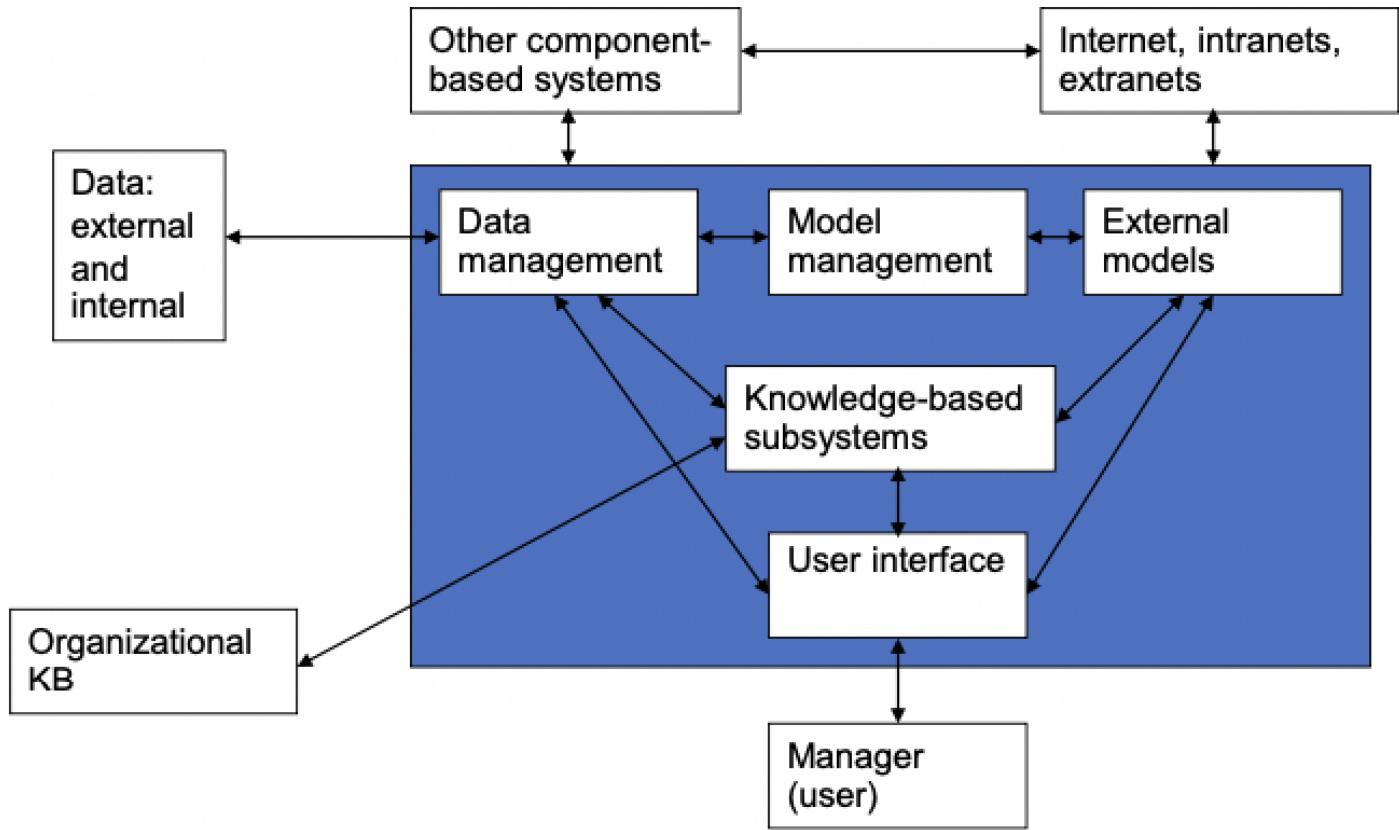
- Key characteristics and capabilities of DSS:



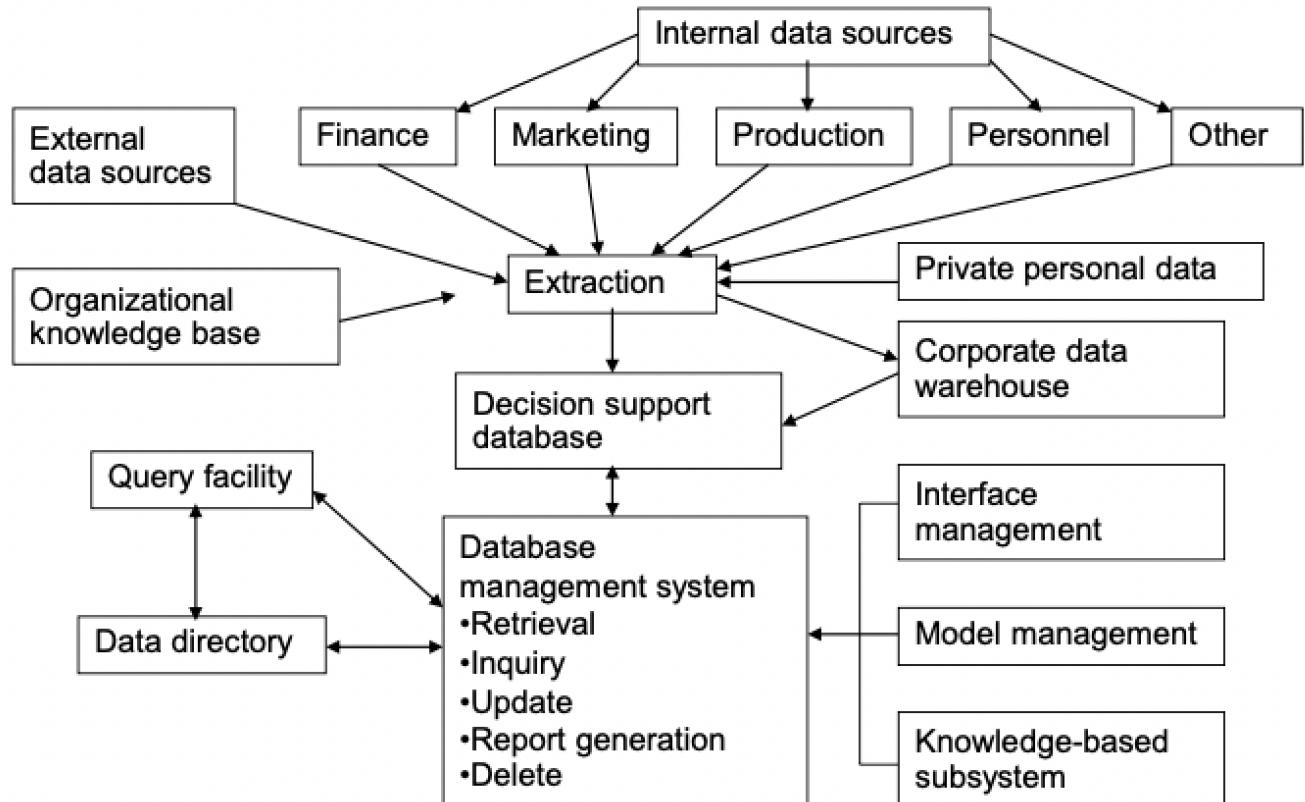
- Components of DSS:

- **Data Management Subsystems**: include Database Management System (DBMS) and data warehouse.
- **Model Management Subsystems**: include financial, statistical, management science, or other quantitative models that provide the analytical capabilities (also called model base management system MDMS).
- **User Interface Subsystems**: include graphical user interface (GUI) that allows users to communicate with the system.

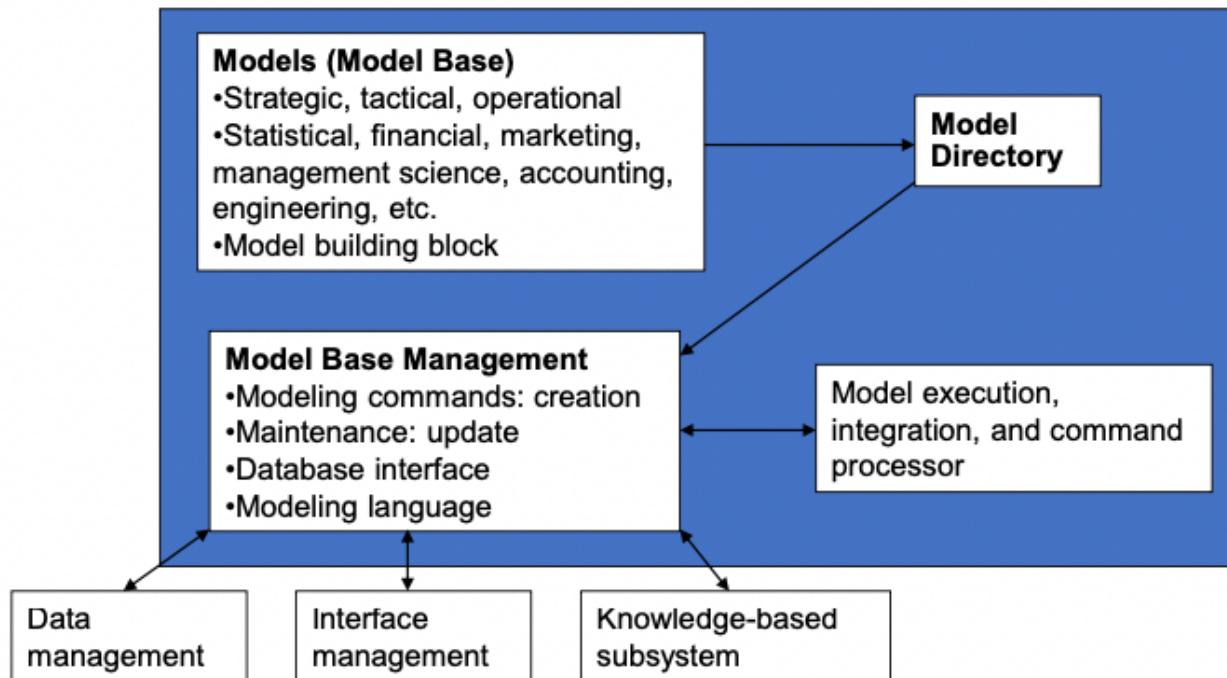
- Schematic View of DSS:



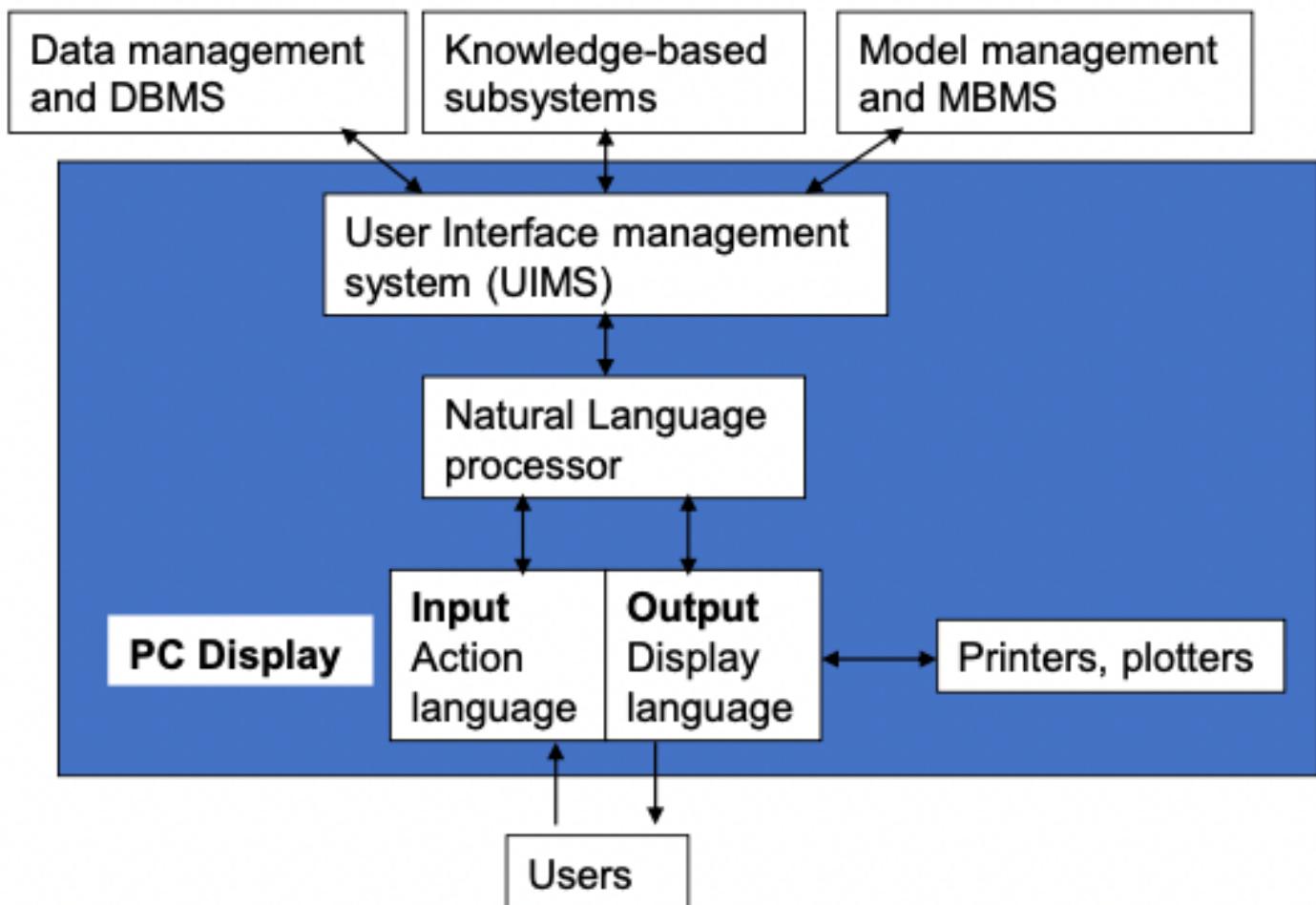
- The Structure of the DMS:



- The Structure of MMS:



- Schematic View of the User Interface Systems:



## Modeling

- Static and Dynamic Models
- Certainty, Uncertainty, and Risk
- Modeling with Spreadsheets
- Decision Tables and Decision Trees
- The Structure of Mathematical Models
- Mathematical Programming Optimization
- Multiple Goals, Sensitivity, What-if and Goal Seeking
- Problem Solving Search Methods

## Static and Dynamic Models

- A static model takes a single snapshot of the system. The decision is made on that snapshot.
- A decision where to buy a product, a quarterly or annual income statement, the decision to invest are static.
- A dynamic model is time dependent.
- Determining how many checkout points should be open in a supermarket – this needs to consider the time of day because different numbers of customers arrive during each hour.

## Certainty, Uncertainty, and Risk

- In decision making under **certainty**, it is assumed that complete knowledge is available so that the decision maker knows exactly what the outcome of each course of action will be.
- Certainty models are relatively easy to develop and solve, and they can yield optimal solutions.
- In decision making under **uncertainty**, the decision maker consider situations in which several outcomes are possible for each course of action.
- The decision maker does not know, or cannot estimate, the probability of occurrence of the possible outcomes.
- It is more difficult than making decision under certainty because there is insufficient information.
- A decision made under risk (also known as a probabilistic, or stochastic) is one in which the decision maker must consider several possible outcomes for each alternative, each with a given probability of occurrence.
- Risk analysis is a decision-making method that analyzes the risk (based on assumed known probabilities) associated with different alternatives.

## Modeling with Spreadsheets

- Spreadsheet packages were quickly recognized as easy-to-use implementation software for the development of a wide range of applications in business, engineering, mathematics, and science.
- Spreadsheets include extensive statistical, forecasting, and other modeling and database management capabilities, functions, and routines.
- These DSS-related spreadsheets include solver (solver.com), What's best (lindo.com), Braincel (promland.com), NeuralTools, Evolver, @RISK (palisade.com), and GRG-2 (MS Excel).

## Decision Tables and Trees

- Decision tables conventionally organize information and knowledge in a systematic tabular manner to prepare for its analysis.
- Decision under certainty:

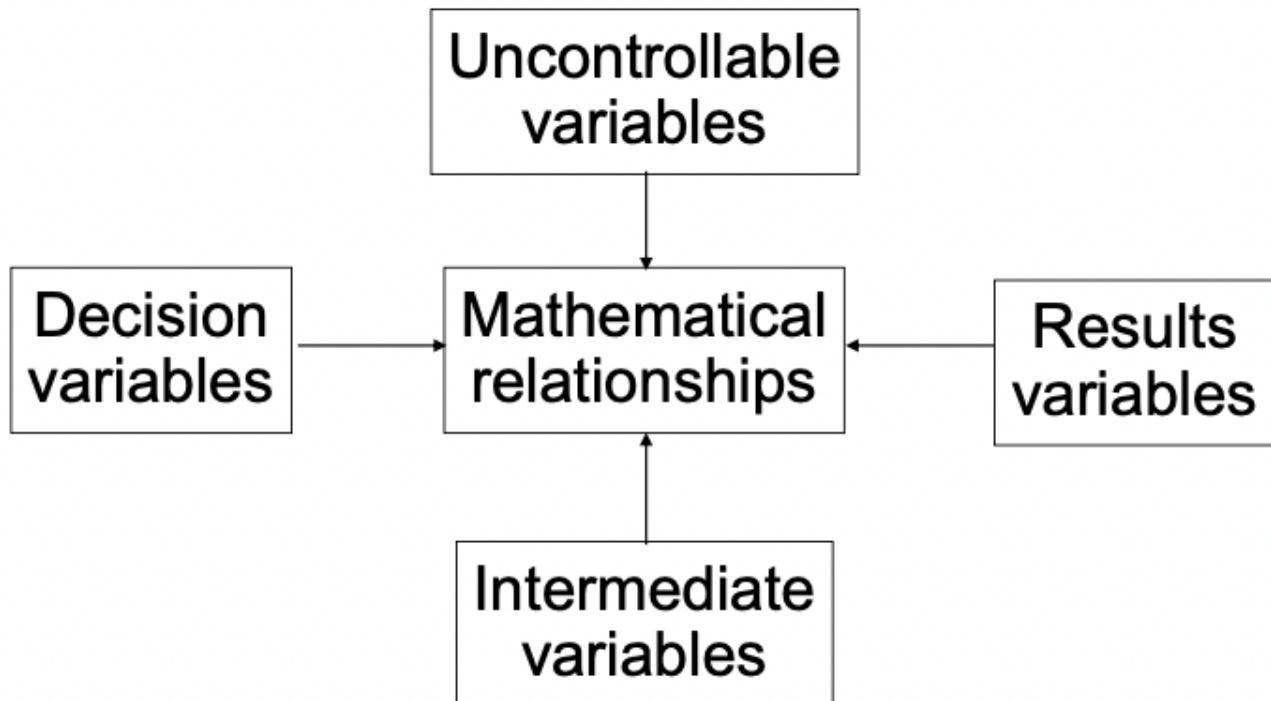
		State of Nature (Uncontrollable Variables)		
Alternative		Solid Growth(%)	Stagnation(%)	Inflation(%)
Bonds		12	6	3
Stocks		15	3	-2
CDs		6.5	6.5	6.5

- Decision under risk:

	Solid Growth	Stagnation	Inflation	Expected Value
Alternative	.5 (%)	.3 (%)	.2 (%)	(%)
Bonds	12	6	3	8.4
Stocks	15	3	-2	8.0
CDs	6.5	6.5	6.5	6.5

- A decision tree shows the relationships of the problem graphically and can handle complex situation in a compact form.
- TreeAge Pro ([treeage.com](http://treeage.com)), PrecisionTree ([palisade.com](http://palisade.com)), [psychwww.com/mtsite/dectree.html](http://psychwww.com/mtsite/dectree.html) and Mind Tools ([mindtools.com](http://mindtools.com)).

## Structure of Mathematical Model



## Example of the Components of Models

Area	Decision variables	Result variables	Uncontrollable variables and parameters
Financial investment	Investment alternatives and amounts	Total profit, risk Rate of return on investment (ROI) Earning per share Liquidity level	Inflation rate Prime rate Competition
Marketing	Advertising budget Where to advertise	Market share Customer satisfaction	Customer's income Competitor's action
Manufacturing	What and how much to produce Inventory levels Compensation programs	Total cost Quantity level Employee satisfaction	Machine capacity Technology Material prices
Accounting	Use of computers Audit schedule	Data processing cost Error rate	Computer technology Tax rates Legal requirements
Transportation	Shipments schedule Use of smart cards	Total transport cost Payment float time	Delivery distance Regulations
Services	Staffing levels	Customer satisfaction	Demand for services

## Mathematical Programming Optimization

- Linear Programming.
  - o Product Mix.
  - o Transportation Problem.
- Non-Linear Programming.
  - o Travelling salesman.
  - o Vehicle routing problem.

## Multiple Goals

- Managers want to attain simultaneous goals, some of which may conflict.
- In addition to earning money, the company wants to grow, develop its products and employees, provide job security to its workers.
- Managers want to satisfy the shareholders and at the same time enjoy high salaries and expense accounts, and employees want to increase their take-home pay and benefits.
- To solve this kind of problems, common methods are:
  - o Utility theory.
  - o Goal programming.
  - o Expression of goals and constraints, using LP.
  - o A point system.

## Sensitivity Analysis

- Sensitivity analysis attempts to assess the impact of a change in the input data or parameters on the proposed solution.
- Sensitivity allows flexibility and adaptation to changing conditions and to the requirements of different decision-making situations.
- Sensitivity analysis tests relationships such as the following:
  - o The impact of changes in external (uncontrollable) variables and parameters on the outcome variables.
  - o The impact of changes in decision variables on the outcome variables.
  - o The effect of uncertainty in estimating external variables.
  - o The effects of different dependent interactions among variables.
  - o The robustness of decisions under changing condition.

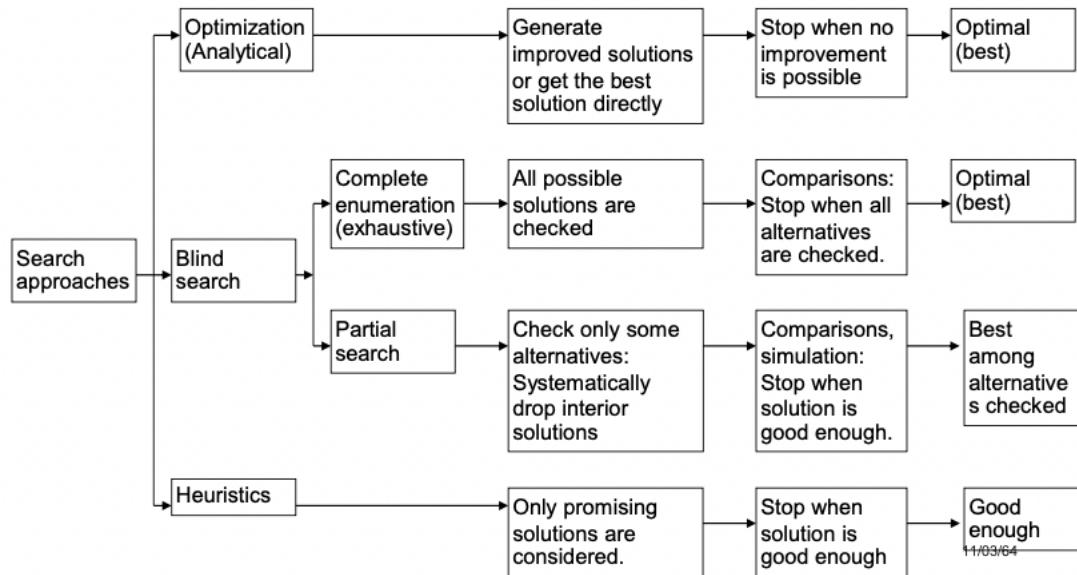
## What-if Analysis

- What-if analysis is structured as What will happen to the solution if an input variable, an assumption, or a parameter value is changed.
- For example, what will happen to the total inventory cost if the cost of the carrying inventories increases by 10 percent?
- A spreadsheet tool is a good example. A manager can analyze a cash flow problem by changing parameter's values and see the differences without any involvement computer programmers.

## Goal Seeking

- Goal seeking calculates the values of the inputs necessary to achieve a desired level of an output (goal). It represents a backward solution approach.
- For example, what annual R&D budget is needed for an annual growth rate of 15 percent by 2012?

## Problem Solving Search Method



## Table of Contents

<i>Decision Making</i> .....	1
<i>Decision Support Systems</i> .....	1
<i>Modeling</i> .....	5
<i>Static and Dynamic Models</i> .....	5
<i>Certainty, Uncertainty, and Risk</i> .....	5
<i>Modeling with Spreadsheets</i> .....	5
<i>Decision Tables and Trees</i> .....	5
<i>Structure of Mathematical Model</i> .....	6
<i>Example of the Components of Models</i> .....	7
<i>Mathematical Programming Optimization</i> .....	7
<i>Multiple Goals</i> .....	7
<i>Sensitivity Analysis</i> .....	8
<i>What-if Analysis</i> .....	8
<i>Goal Seeking</i> .....	8
<i>Problem Solving Search Method</i> .....	8

# Lecture 6.2: Introduction to Optimization

## Learning Objectives

- Recognize decision-making situations which may benefit from an optimization modeling approach.
- Formulate algebraic models for linear programming problems.
- Develop spreadsheet models for linear programming problems.
- Use Excel's Solver Add-In to solve linear programming problems.
- Interpret the results of models and perform basic sensitivity analysis.

## Optimization: Basic Ideas

- Major field within the discipline of Data Analytics, Operations Research and Management Science.
- Optimization Problem Components:
  - o Decision Variables.
  - o Objective Function (to maximize or minimize).
  - o Constraints (requirements or limitations).
- Basic Idea:
  - o Find the values of the decision variables that maximize (minimize) the objective function value, while staying within the constraints.

## Linear Programming (LP)

- If the objective function and all constraints are linear functions of the decision variables (e.g., no squared terms, trigonometric functions, ratios of variables), then the problem is called a Linear Programming (LP) problem. LPs are much easier to solve by computer than problems involving nonlinear functions.
- Real-world LPs are solved which contain hundreds of thousands to millions of variables (with specialized software!). Our problems are obviously much smaller, but the basic concepts are much the same.

## Real-World Examples

- Dynamic and Customized Pricing.
- Product Mix.
- Scheduling/Allocation.
- Routing/Logistics.
- Supply Chain Optimization.
- Facility Location.
- Financial Planning/Asset Management.
- Etc.

## Solving Optimization Problems

- Understand the problem; draw a diagram.
- Write a problem formulation in words:
- Write the algebraic formulation of the problem:
  - o Define the decision variables.
  - o Write the objective function in terms of the decision variables.
  - o Write the constraints in terms of the decision variables.
- Develop Spreadsheet Model.
- Setup the Solver settings and solve the problem.
- Examine results, make corrections to model and/or Solver settings.

- Interpret the results and draw insights.

## Solver

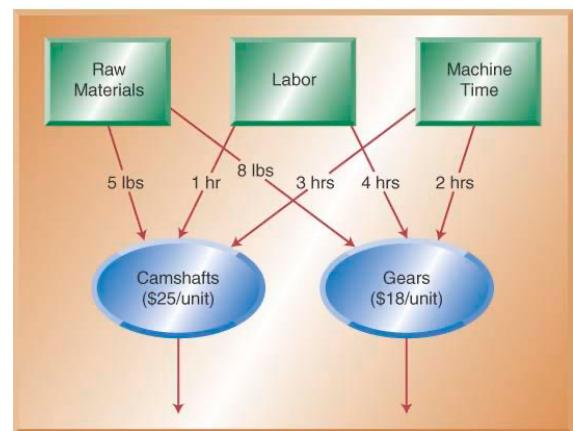
- We will use Solver as simple tool to illustrate the example, an Excel add-in, to solve Linear Programming problems.
- Add-In: An additional piece of software that Excel loads into memory when needed.

## Example: Product Mix Decision

- DJJ Enterprises makes automotive parts, Camshafts & Gears.
- Unit Profit: Camshafts \$25/unit, Gears \$18/unit.
- Resources needed: Steel, Labor, Machine Time. In total, 5000 lbs. steel available, 1500 hours labor, and 1000 hours machine time.
- Camshafts need 5 lbs. steel, 1 hour labor, 3 hours machine time.
- Gears need 8 lbs. steel, 4 hours labor, 2 hours machine time.
- How many camshafts & gears to make to maximize profit?

## Understanding the Problem

- Text-Based Formulation:
  - o Decision Variables: Number of camshafts to make, number of gears to make.
  - o Objective Function: Maximize profit.
  - o Constraints: Don't exceed amounts available of steel, labor, and machine time.



## Algebraic Formulation

- Decision Variables:
  - o  $C$  = number of camshafts to make.
  - o  $G$  = number of gears to make.
- Objective Function:
  - o Maximize  $25C + 18G$  (profit in \$).
- Constraints:
  - o  $5C + 8G \leq 5000$  (steel in lbs.)
  - o  $1C + 4G \leq 1500$  (labor in hours)
  - o  $3C + 2G \leq 1000$  (machine time in hours)
  - o  $C \geq 0, G \geq 0$  (non-negativity)

## Important Concepts

- Linear Program: The objective function and constraint are linear functions of the decision variables. Therefore, this is a Linear Program.
- Feasibility:
  - o Feasible Solution. A solution is feasible for an LP if all constraints are satisfied.
  - o Infeasible Solution. A solution is infeasible if one or more constraints is violated.
  - o Check the solutions  $C=75, G=200$ ; and  $C=300, G=200$  for feasibility.
- Optimal Solution:
  - o The optimal solution is the feasible solution with the largest (for a max problem) objective value (smallest for a min problem).

## Solving Linear Programming Problems

- Trial and error:
  - o Possible for very small problems.
  - o Virtually impossible for large problems.
- Graphical approach:
  - o It is possible to solve a 2-variable problem graphically to find the optimal solution (not shown).
- Simplex Method:
  - o This is a mathematical approach developed by George Dantzig.
  - o Can solve small problems by hand.
- Computer Software:
  - o Most optimization software uses the Simplex Method to solve the problems.
  - o Excel's Solver Add-In is an example of such software.
- Solver can solve LPs of up to 200 variables.

## Spreadsheet Model Development

- Develop correct, flexible, documented model.
- Sections for decision variables, objective function, and constraints.
- Use algebraic formulation and natural structure of the problem to guide structure of the spreadsheet.
- Use one cell for each decision variable.
- Store objective function coefficients in separate cells and use another cell to compute the objective function value.
- Store constraint coefficients in cells, compute the LHS value of each constraint, for comparison to the RHS value.

## Spreadsheet Model

- C=75, G=200 (cells B5, C5) entered as trial values. This is a feasible, but not the optimal, solution.
- Note close relationship to algebraic formulation.
- Note that only one distinct formula needed to be entered; once entered in Cell D8, it was copied to Cells D11:D13.
- This is possible because the coefficients were stored separately in a specific structure.

	A	B	C	D	E	F	G
1	<b>Example B.1</b>						
2	<b>DJJ Enterprises Production Planning</b>						
3							
4	<b>Decision Variables</b>	Camshafts	Gears				
5	Units to Make	75	200				
6							
7	<b>Objective</b>			<b>Total</b>			
8	Profit	\$25	\$18	\$5,475	D8: =B8*B\$5+C8*C\$5 (copied to D11:D13)		
9							
10	<b>Constraints</b>			<b>Used</b>		<b>Available</b>	
11	Steel (lbs)	5	8	1975	<=	5000	
12	Labor (hrs)	1	4	875	<=	1500	
13	Machine Time (hrs)	3	2	625	<=	1000	

## Solver Basics

- Don't even think about using Solver until you have a working, flexible spreadsheet model that you can use as a "what if" tool!
- Solver Settings
  - o Specify Objective Cell (objective function).
  - o Specify Changing Cells (decision variables).
  - o Specify Constraints.
  - o Specify Solver Settings
- Solve Problem to find Optimal Solution.

## Solver Settings: Target Cell and Changing Cells

- Specify Target Cell: D8.
  - o Equal to: Max.
- Changing Cells (decision variables)
  - o B5, C5.

## Solver Settings: Constraints

- Click "Add" to add constraints.
- Select LHS Cell (D11), relationship ( $\leq$ ), and RHS Cell (F11).
- LHS Cell should contain a formula which computes the LHS Value of the constraint.
- Typically, RHS Cell should contain a fixed value, but this is not absolutely required.
- Repeat for the other two constraints (for labor and machine time).
- **OR use the range for the LHS and the RHS** (see next slide).

## Solver Options

- Select Simplex LP"
- Tells Solver to use the Simplex Method, which is faster and is Solver's default optimization method.
- Check "Make Unconstrained Variables Non-negative".
- Tells Solver that the decision variables (B5:C5, representing the number of Camshafts and Gears) must be  $\geq 0$  in any feasible solution.
- Leave other settings at their defaults.

## Completed Solver Box

- Click "Solve" to tell Solver to find the Optimal Solution.

## Solved Spreadsheet (Optimal Solution)

- Optimal Solution: Make 100 camshafts, 350 gears.
- Optimal Objective Value: \$8800 profit.
- Both pieces of information are important. Knowing the optimal objective value is useless without knowing how that value can be attained.

	A	B	C	D	E	F
1	<b>Example B.1</b>					
2	<b>DJJ Enterprises Production Planning</b>					
3						
4	<b>Decision Variables</b>	Camshafts	Gears			
5	Units to Make	100	350			
6						
7	<b>Objective</b>			Total		
8	Profit	\$25	\$18	\$8,800		
9						
10	<b>Constraints</b>			Used		Available
11	Steel (lbs)	5	8	3300	$\leq$	5000
12	Labor (hrs)	1	4	1500	$\leq$	1500
13	Machine Time (hrs)	3	2	1000	$\leq$	1000

## Interpreting the Solution

- Which constraints are actively holding us back from making even more profit?
  - o Binding Constraints: Constraints at the optimal solution with the LHS equal to the RHS; equivalently for a resource, a binding constraint is one in which all the resource is used up.
  - o Labor (all 1500 hours used), Machine Time (all 1000 hours used).
- Non-Binding Constraints:
  - o Steel: Have 5000 lbs. available, but only need 3300 lbs. for this solution.
- What-If Analysis. Once the Spreadsheet and Solver model is set up, it is easy to change one or more input values and re-solve the problem.
  - o This interactive use can be a very powerful way to use optimization.

## Solution Reports

- Solver can generate three solution reports:
  - o Answer Report.
  - o Sensitivity Report.
  - o Limits Report: Not covered here.
- The Answer Report presents in a standard format the Solver Settings and the optimal solution.
- The Sensitivity Report shows what will happen if certain problem parameters are changed from their current values.

## Answer Report

- Tightening a binding constraint can only worsen the objective function value, and loosening a binding constraint can only improve the objective function value. As such, once an optimal solution is found, managers can seek to improve that solution by finding ways to relax binding constraints.
- Three sections: Objective Cell, Changing Cells, Constraints.
- “Final Value” indicates optimal solution.
- For each constraint, binding/non-binding status, and “slack” (difference between LHS and RHS), slack is zero for binding constraints.
- Also note that the Solver Settings for Objective Cell, Changing Cells, and Constraints are reported here. This can be a useful debugging tool.

Objective Cell (Max)					
Cell	Name	Original Value	Final Value		
\$D\$8	Profit Total	\$0	\$8,800		
Variable Cells					
Cell	Name	Original Value	Final Value	Integer	
\$B\$5	Units to Make Camshafts	0	100	Contin	
\$C\$5	Units to Make Gears	0	350	Contin	
Constraints					
Cell	Name	Cell Value	Formula	Status	Slack
\$D\$11	Steel (lbs) Used	3300	\$D\$11<=\$F\$11	Not Binding	1700
\$D\$12	Labor (hrs) Used	1500	\$D\$12<=\$F\$12	Binding	0
\$D\$13	Machine Time (hrs) Used	1000	\$D\$13<=\$F\$13	Binding	0

## Sensitivity Report

Variable Cells						
Cell	Name	Final Value	Reduced Cost	Objective Coefficient	Allowable Increase	Allowable Decrease
\$B\$5	Units to Make Camshafts	100	0	25	2	20.5
\$C\$5	Units to Make Gears	350	0	18	82	1.333333333
Constraints						
Cell	Name	Final Value	Shadow Price	Constraint R.H. Side	Allowable Increase	Allowable Decrease
\$D\$11	Steel (lbs) Used	3300	0	5000	1E+30	1700
\$D\$12	Labor (hrs) Used	1500	0.4	1500	500	1166.666667
\$D\$13	Machine Time (hrs) Used	1000	8.2	1000	1416.666667	250

- Note two sections:
  - o Variable Cells.
  - o Constraints.
- **Constraint Section**
  - o Shadow Price: This is the amount the optimal objective value will change by if the RHS of the constraint is increased by one unit.
  - o Units of shadow price: (objective function units/constraint units). Example:
    - For the steel constraint, the units are \$/lb.
    - For the labor constraint, the units are \$/hour.

- Allowable Increase/Decrease: Provides the increase/decrease of the RHS of the constraint for which the shadow price stays the same; that is, the effect on the objective value stays the same in this range.
- Example:
  - Shadow Price of Machine Time is \$8.20/hour, valid for an increase of 1416 hours or a decrease of 250 hours.
  - If we make 200 additional machine time hours available, profit can be increased by (200 hours) (\$8.20/hour) = \$1640.
- What is the shadow price of a non-binding constraint? Why? Will this always be the case?
- **Variable Cells Section**
- Allowable Increase/Decrease: This is the amount the objective coefficient for a decision variable can be increased/decreased without changing the optimal solution.
- Example
  - Gears currently have a profit of \$18/unit (objective coefficient).
  - Allowable Increase/Decrease is \$82 and \$1.33, respectively.
  - Interpretation: If the unit profit for gears is between \$16.67 and \$100, the optimal solution (production plan) will be to make 100 camshafts and 350 gears.
  - Note: The optimal solution (production plan) stays the same within this range, but the optimal objective value (profit) changes since the unit profit is changing. So, if the profit on gears goes up to \$20/unit, profit will increase by (\$20-\$18) (350) = \$700.

## Sensitivity Report: Caveat

- Effects identified in the Sensitivity Report need to be carefully interpreted.
- Specifically, the effect of the change of one parameter (e.g., a RHS value or an objective coefficient) assume that all other parameters of the model stay at their base case values.
- For example, the Sensitivity Report does not tell us what happens if additional quantities of both Labor and Machine Time become available. In this scenario, we would need to enter the new values and re-solve the model.

## Highlights

- People use informal “optimization” to make decisions almost every day.
- Organizations use formal optimization methods to address problems across the organization, from optimal pricing to locating a new facility.
- The algebraic formulation of an LP comprises the definitions of the decision variables, an algebraic statement of the objective function, and algebraic statements of the constraints.
- The spreadsheet model for an optimization problem should be guided by the algebraic formulation.
- Solver, an Excel Add -In, can solve both linear and nonlinear problems. This lecture focuses on solving linear problems.
- After solving an LP, you must interpret the results to see if they make sense, fix problems with the model, and find the insights useful for management.
- Solver can generate the Answer and Sensitivity Reports. The Sensitivity Report provides additional information about what happens to the solution when certain coefficients of the problem are changed.

## Table of Contents

<i>Learning Objectives</i> .....	1
<i>Optimization: Basic Ideas</i> .....	1
<i>Linear Programming (LP)</i> .....	1
<i>Real-World Examples</i> .....	1
<i>Solving Optimization Problems</i> .....	1
<i>Solver</i> .....	2
<i>Example: Product Mix Decision</i> .....	2
<i>Understanding the Problem</i> .....	2
<i>Algebraic Formulation</i> .....	2
<i>Important Concepts</i> .....	2
<i>Solving Linear Programming Problems</i> .....	3
<i>Spreadsheet Model Development</i> .....	3
<i>Spreadsheet Model</i> .....	3
<i>Solver Basics</i> .....	3
<i>Solver Settings: Target Cell and Changing Cells</i> .....	4
<i>Solver Settings: Constraints</i> .....	4
<i>Solver Options</i> .....	4
<i>Completed Solver Box</i> .....	4
<i>Solved Spreadsheet (Optimal Solution)</i> .....	4
<i>Interpreting the Solution</i> .....	4
<i>Solution Reports</i> .....	5
<i>Answer Report</i> .....	5
<i>Sensitivity Report</i> .....	5
<i>Sensitivity Report: Caveat</i> .....	6
<i>Highlights</i> .....	6

# Lecture 7: Introduction to Data Mining

## Introduction

- Data is growing at a phenomenal rate.
- Users expect more sophisticated information.
- How?

UNCOVER HIDDEN INFORMATION  
**DATA MINING**

## Data Mining – Definition

- Finding hidden information in a database.
- Fit data to a model.
- Similar terms:
  - o Exploratory data analysis.
  - o Data driven discovery.
  - o Deductive learning.

## Data Mining Algorithm

- Objective: Fit Data to a Model:
  - o Descriptive.
  - o Predictive.
- Preference – Technique to choose the best model.
- Search – Technique to search the data.
  - o “Query”.

## Database Processing vs. Data Mining Processing

<ul style="list-style-type: none"><li>• <b>Query</b><ul style="list-style-type: none"><li>• Well defined</li><li>• SQL</li></ul></li><li>■ <b>Data</b><ul style="list-style-type: none"><li>– Operational data</li></ul></li><li>■ <b>Output</b><ul style="list-style-type: none"><li>– Precise</li><li>– Subset of database</li></ul></li></ul>	<ul style="list-style-type: none"><li>• <b>Query</b><ul style="list-style-type: none"><li>• Poorly defined</li><li>• No precise query language</li></ul></li><li>■ <b>Data</b><ul style="list-style-type: none"><li>– Not operational data</li></ul></li><li>■ <b>Output</b><ul style="list-style-type: none"><li>– Fuzzy</li><li>– Not a subset of database</li></ul></li></ul>
--	--

## Query Examples

- Database:
  - o Find all credit applicants with last name of Smith.
  - o Identify customers who have purchased more than \$10.000 in the last month.
- Data Mining:
  - o Find all credit applicants who are poor credit risks (classification).
  - o Identify customers with similar buying habits (clustering).
  - o Find all items which are frequently purchased with milk (association rules).

## Basic Data Mining Tasks

- Classification maps data into predefined groups or classes.
  - o Supervised learning.
  - o Prediction.
  - o Regression.
- Clustering groups similar data together into clusters.
  - o Unsupervised learning.
  - o Segmentation.
  - o Partitioning.
- Link Analysis uncovers relationships among data.
  - o Affinity Analysis.
  - o Association Rules.
  - o Sequential Analysis determines sequential patterns.

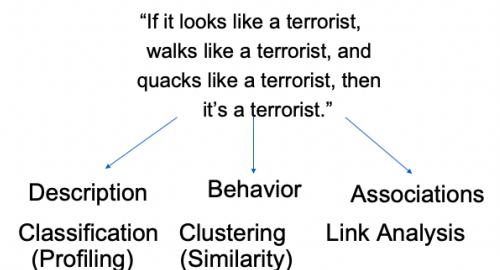
## Classification

- Assign data into predefined groups or classes.

## But It Isn't Magic

- You must know what you are looking for.
- You must know how to look for you.
- Suppose you know that a specific cave had gold:
  - o What would you look for?
  - o How would you look for it?
  - o Might need an expert miner.

"If it looks like a duck,  
walks like a duck, and  
quacks like a duck, then  
it's a duck."



## Classification Example

- Grading.
- Given a collection of annotated data (in this case 5 instances of Katydids and five of Grasshoppers), decide what type of insect the unlabeled example is.
- The classification problem can now be expressed as:
  - o Given a training database, predict the class label of previously unseen instance.
- Facial Recognition.
- Handwriting Recognition.
- Anomaly Detection.

## Clustering

- Partition data into previously undefined groups.

## Two Types of Clustering

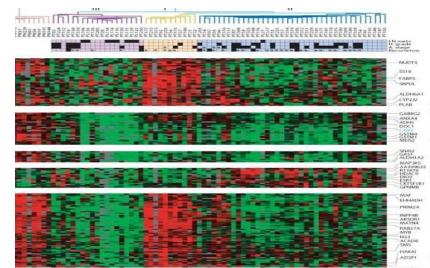
- Hierarchical.
- Partitional.

## Hierarchical Clustering Example

- Iris Data Set.

## Microarray Data Analysis

- Each probe location associated with gene.
- Color indicates degree of gene expression.
- Compare different samples (normal/disease).
- Track same sample over time.
- Questions:
  - o Which genes are related to this disease?
  - o Which genes behave in a similar manner?
  - o What is the function of a gene?
- Clustering:
  - o Hierarchical.
  - o K-means.
- Gene Expression Profiling identifies clinically relevant subtypes of prostate cancer.



## Association Rules/Link Analysis

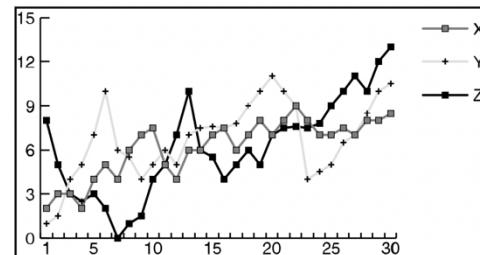
- Find relationships between data.

## Association Rules Examples

- People who buy diapers also buy beer.
- If gene A is highly expressed in this disease, then gene A is also expressed.
- Relationships between people.
- Book Stores.
- Department Stores.
- Advertising.
- Product Placement.

## Example: Stock Market Analysis

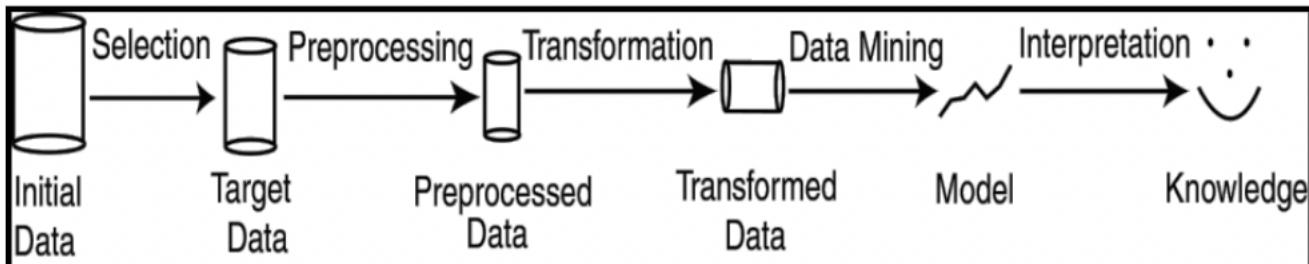
- Example: Stock Market.
- Predict future values.
- Determine similar patterns over time.
- Classify behavior.



## Data Mining vs. KDD

- Knowledge Discovery in Databases (KDD): process of finding useful information and patterns in data.
- Data Mining: Use of algorithms to extract the information and patterns derived by the KDD process.

## KDD Process



- Selection: Obtain data from various sources.
- Preprocessing: Cleanse data.
- Transformation: Convert to common format. Transform to new format.
- Data Mining: Obtain desired results.
- Interpretation/Evaluation: Present results to user in meaningful manner.

## KDD Process Example: Web Log

- Selection:
  - o Select log data (dates and locations) to use.
- Preprocessing:
  - o Remove identifying URLs. Remove error logs.
- Transformation:
  - o Sectionize (Sort and Group).
- Data Mining:
  - o Identify and count patterns. Construct data structure.
- Interpretation/Evaluation:
  - o Identify and display frequently accessed sequences.
- Potential User Applications:
  - o Cache prediction.
  - o Personalization.

## Related Topics

- Databases.
- OLTP.
- OLAP.
- Information Retrieval.

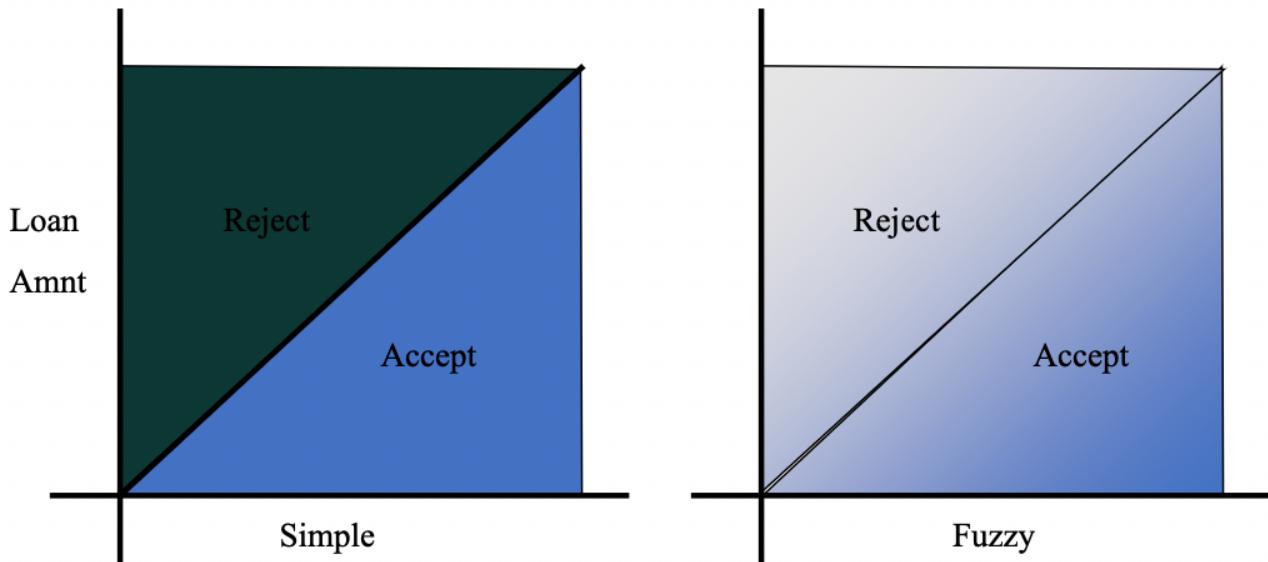
## DB & OLTP Systems

- Schema:
  - o (ID, Name, Address, Salary, JobNo).
- Data Model:
  - o Entity Relationship.
  - o Relational.
- Transaction.
- Sample query:

```
SELECT Name
FROM T
WHERE Salary > 100000
```

***DM: Only imprecise queries***

## Classification/Prediction is Fuzzy



## Information Retrieval

- **Information Retrieval (IR)**: retrieving desired information from textual data.
- Library Science.
- Digital Libraries.
- Web Search Engines.
- Traditionally keyword based.
- Sample query:
  - o Find all documents about “data mining”.

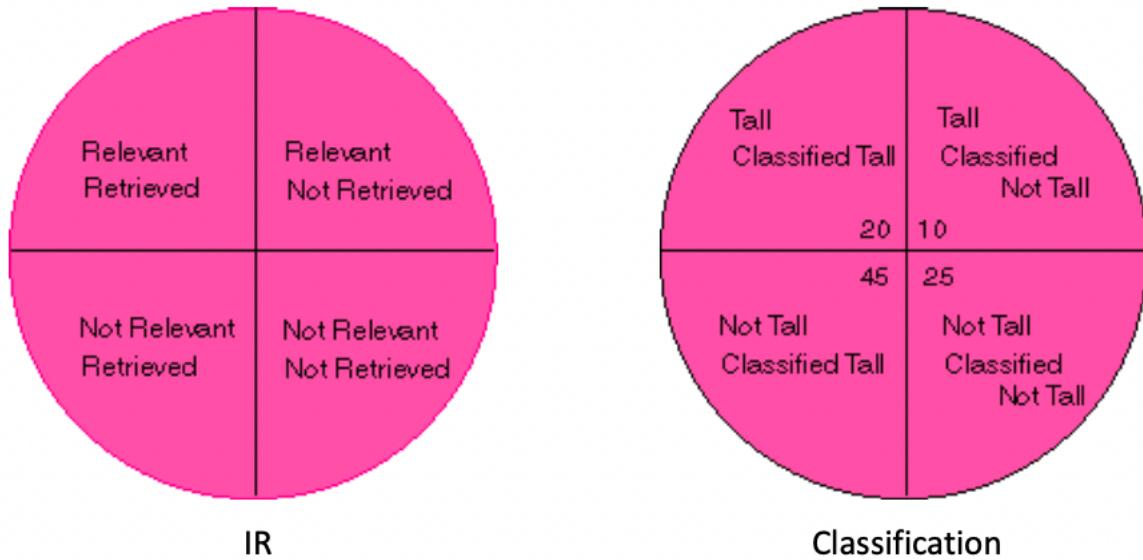
## *DM: Similarity measures;*

*Mine text/Web data.*

- **Similarity**: measure of how close a query is to a document.
- Documents which are “close enough” are retrieved.
- Metrics:

- **Precision** =  $\frac{|\text{Relevant and Retrieved}|}{|\text{Retrieved}|}$
- **Recall** =  $\frac{|\text{Relevant and Retrieved}|}{|\text{Relevant}|}$

## IR Query Result Measures and Classification



## OLAP

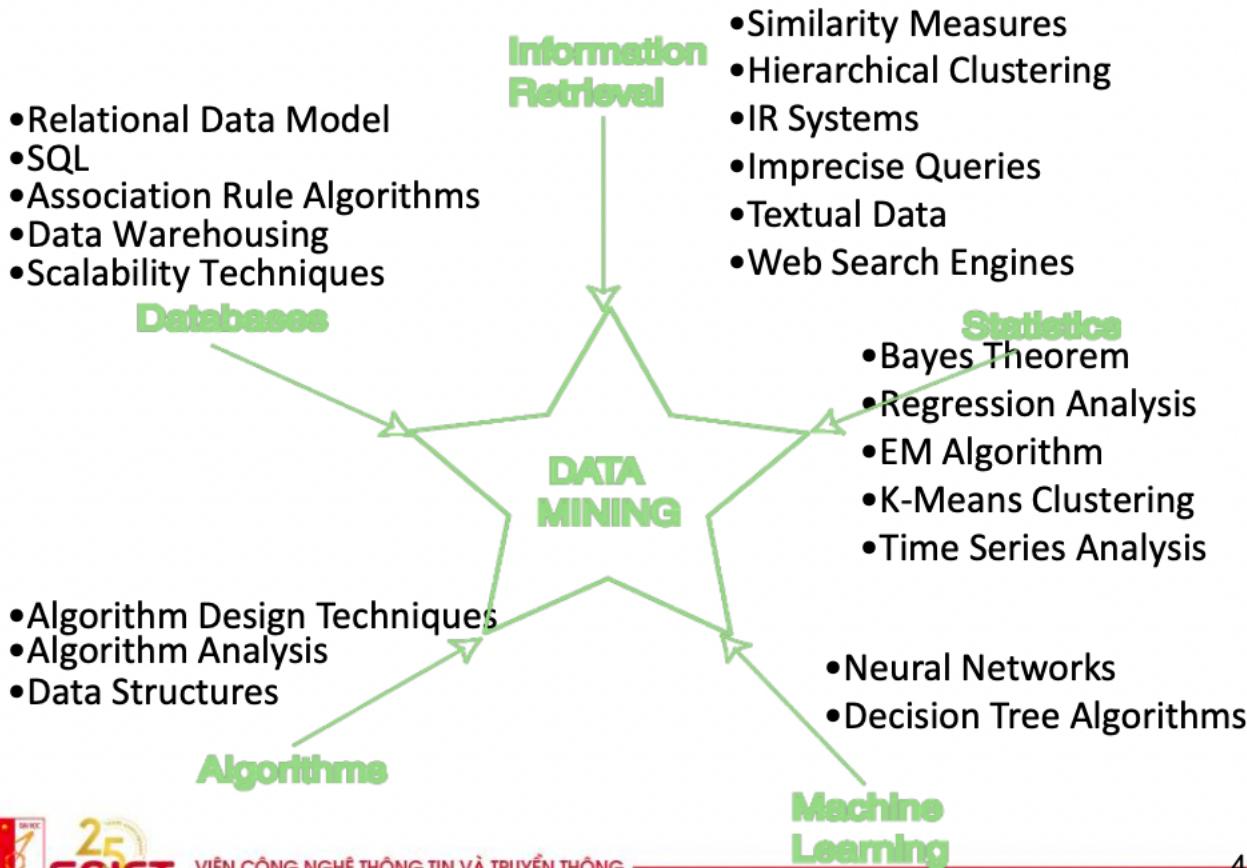
- **Online Analytic Processing (OLAP)**: provides more complex queries than OLTP.
- **Online Transaction Processing (OLTP)**: traditional database/transaction processing.
- Dimensional data; cube view.
- Visualization of operations:
  - o **Slice**: examine sub-cube.
  - o **Dice**: rotate cube to look at another dimension.
  - o **Roll Up/Drill Down**.

***DM: May use OLAP queries.***

## DM vs. Related Topics

Area	Query	Data	Results	Output
DB/OLTP	Precise	Database	Precise	DB Objects or Aggregation
IR	Precise	Documents	Vague	Documents
OLAP	Analysis	Multidimensional	Precise	DB Objects or Aggregation
DM	Vague	Preprocessed	Vague	KDD Objects

# Data Mining Development



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

49

## KDD Issues

- |   |  |
|---|--|
| <ul style="list-style-type: none"><li>- Human Interaction.</li><li>- Overfitting.</li><li>- Outliers.</li><li>- Interpretation.</li><li>- Visualization.</li><li>- Large Datasets.</li><li>- High Dimensionality.</li></ul> | <ul style="list-style-type: none"><li>- Multimedia Data.</li><li>- Irrelevant Data.</li><li>- Noisy Data.</li><li>- Changing Data.</li><li>- Integration.</li><li>- Application.</li></ul> |
|---|--|

## Warning

- With data mining, you don't always know what you are looking for.
- There is not one right answer.
- The data you are using is noisy.
- Data Mining is a very applied discipline.
- A data mining course provides you tools to use to analyze data.
- Experience provides you knowledge of how to use these tools.

## Social Implications of DM

- Privacy. - Profiling.
- Unauthorized use. - Invalid results and claims.

## Data Mining Metrics

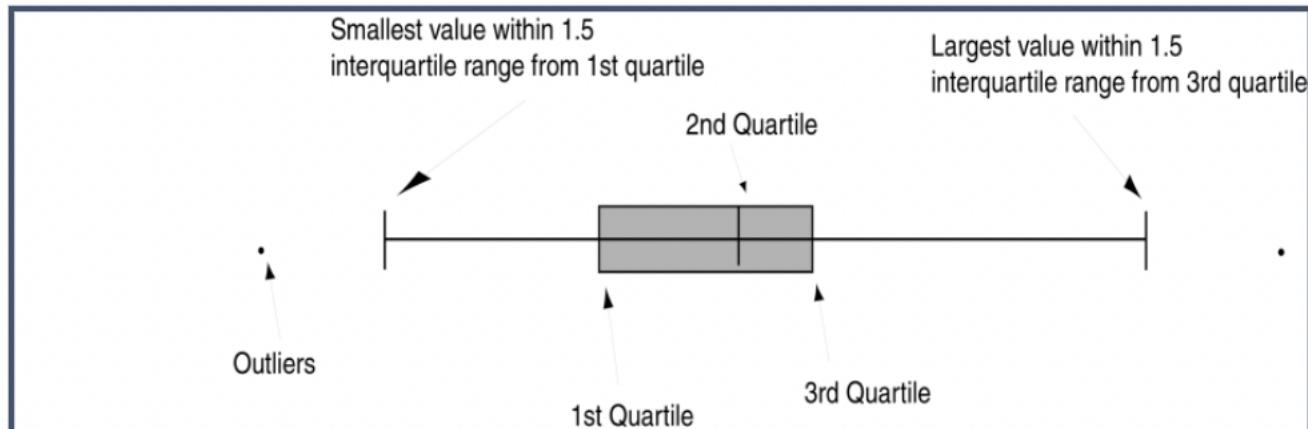
- Usefulness.
- Return on Investment (ROI).
- Accuracy.
- ...
- Space/Time.

## Visualization Techniques

- Graphical.
- Geometric.
- Icon-based.
- Pixel-based.
- Hierarchical.
- Hybrid.

## Models Based on Summarization

- Visualization: Frequency distribution, mean, variance, median, mode, etc.
- Box Plot:



## DM Tools

- XLMiner – Easy addin to Excel: <http://www.solver.com/xlminer/index.html>
- Webka – Open Source; Visualization, Functionality, Interface: <http://www.cs.waikato.ac.nz/ml/weka/>
- SAS (JMP) – Commercial Product.
- SPSS – Commercial Product.
- MATLAB – Statistical/Math Applications.
- R – Programming.

## Table of Contents

<i>Introduction</i> .....	1
<i>Data Mining – Definition</i> .....	1
<i>Data Mining Algorithm</i> .....	1
<i>Database Processing vs. Data Mining Processing</i> .....	1
<i>Query Examples</i> .....	1
<i>Basic Data Mining Tasks</i> .....	2
<i>Classification</i> .....	2
<i>But It Isn't Magic</i> .....	2
<i>Classification Example</i> .....	2
<i>Clustering</i> .....	2
<i>Two Types of Clustering</i> .....	2
<i>Hierarchical Clustering Example</i> .....	2
<i>Microarray Data Analysis</i> .....	3
<i>Association Rules/Link Analysis</i> .....	3
<i>Association Rules Examples</i> .....	3
<i>Example: Stock Market Analysis</i> .....	3
<i>Data Mining vs. KDD</i> .....	3
<i>KDD Process</i> .....	3
<i>KDD Process Example: Web Log</i> .....	4
<i>Related Topics</i> .....	4
<i>DB &amp; OLTP Systems</i> .....	4
<i>Classification/Prediction is Fuzzy</i> .....	5
<i>Information Retrieval</i> .....	5
<i>IR Query Result Measures and Classification</i> .....	6
<i>OLAP</i> .....	6
<i>DM vs. Related Topics</i> .....	6
<i>Data Mining Development</i> .....	7
<i>KDD Issues</i> .....	7
<i>Warning</i> .....	7
<i>Social Implications of DM</i> .....	7
<i>Data Mining Metrics</i> .....	8
<i>Visualization Techniques</i> .....	8
<i>Models Based on Summarization</i> .....	8
<i>DM Tools</i> .....	8

# Lecture 8: Classification

## Classification: Definition

- Given a collection of records (training set).
  - o Each record is characterized by a tuple  $(x, y)$  where  $x$  is the attribute set and  $y$  is the class label.
- Task: Learn a model that maps each attribute set  $x$  into one of the predefined class labels  $y$ .
- Example:

Task	Attribute set, $x$	Class label, $y$
Categorizing email messages	Features extracted from email message header and content	spam or non-spam
Identifying tumor cells	Features extracted from x-rays or MRI scans	malignant or benign cells
Cataloging galaxies	Features extracted from telescope images	Elliptical, spiral, or irregular-shaped galaxies

## General Approach for Building Classification Model

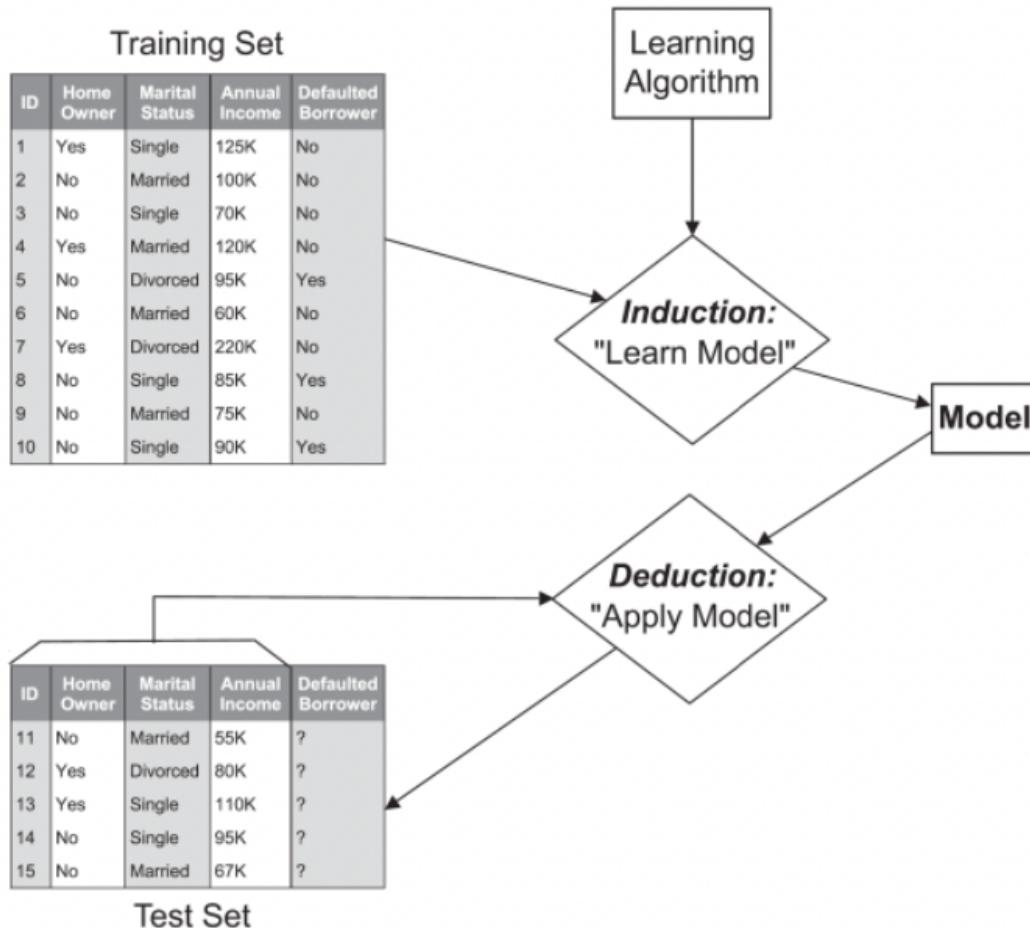
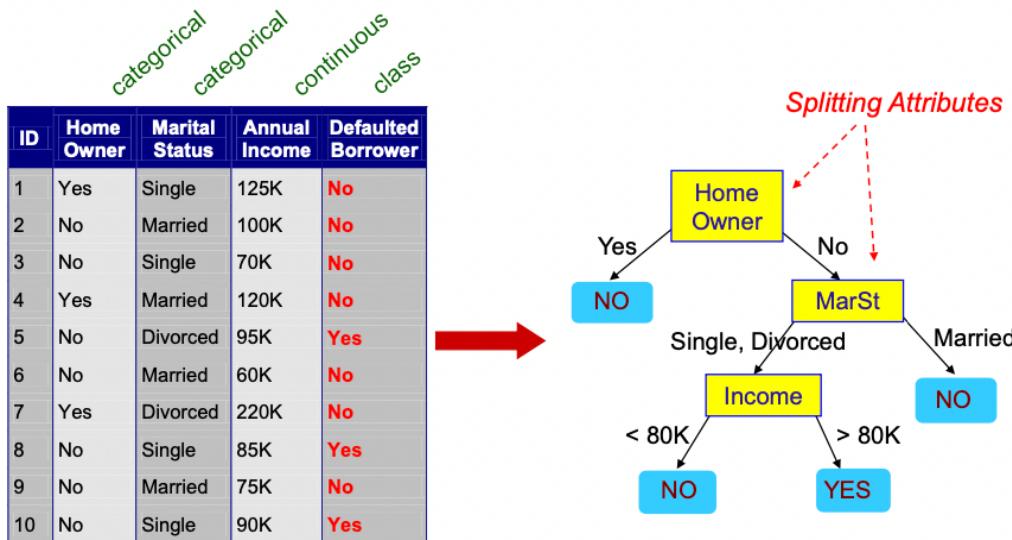


Figure 3.3. General framework for building a classification model.

## Classification Techniques

- Base Classifiers:
  - o Decision Tree based Methods.
  - o Rule-based Methods.
  - o Nearest-neighbor.
  - o Naïve Bayes and Bayesian Belief Networks.
  - o Support Vector Machines.
- Ensemble Classifiers:
  - o Booting.
  - o Bagging.
  - o Random Forests.

### Example of a Decision Tree

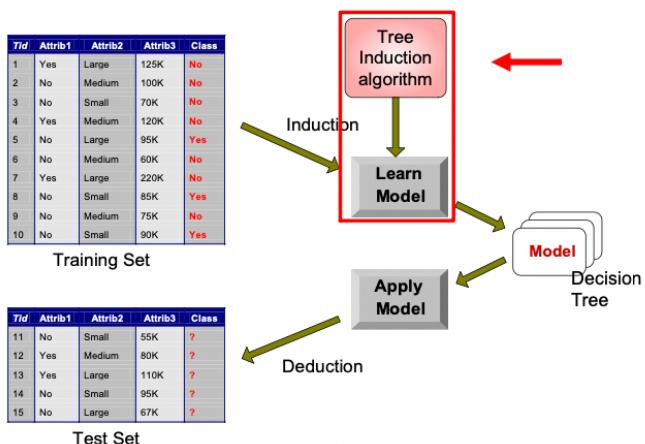


Training Data

Model: Decision Tree

- Apply Model to Test Data?
- There could be more than one tree that fits the same data.

## Decision Tree Classification Task



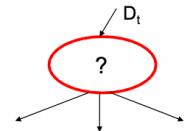
## Decision Tree Induction

- Many Algorithms:
  - o Hunt's Algorithm (one of the earliest).
  - o CART.
  - o ID3, C4.5.
  - o SLIQ, SPRINT.

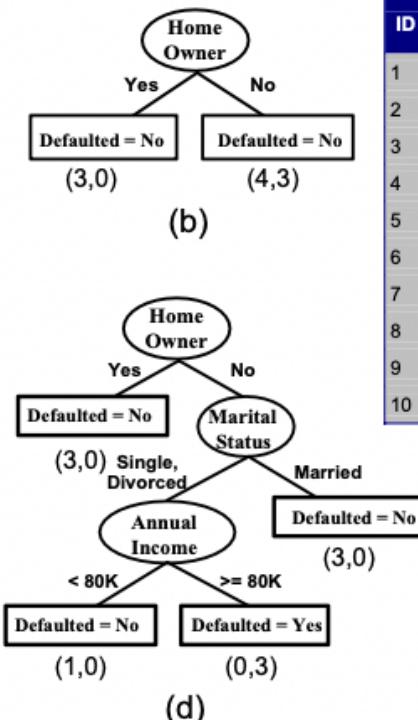
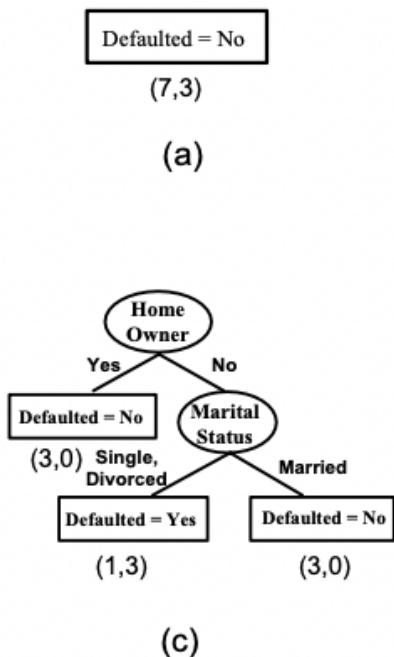
## General Structure of Hunt's Algorithm

- Let  $D_t$  be the set of training records that reach a node  $t$ .
- General Procedure:
  - o If  $D_t$  contains records that belong to the same class  $y_t$  then  $t$  is a leaf node labeled as  $y_t$ .
  - o If  $D_t$  contains records that belong to more than one class, use an attribute test to split the data into smaller subsets.
  - o Recursively apply the procedure to each subset.

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
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



## Hunt's Algorithm



ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
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

## Design Issues of Decision Tree Induction

- How should training records be split?
  - o Method for expressing test condition.
    - Depending on attribute types.
  - o Measure for evaluating the goodness of a test condition.
- How should the splitting procedure stop?
  - o Stop splitting if all the records belong to the same class or have identical attribute values.
  - o Early termination.

## Methods for Expressing Test Conditions

- Depends on attribute types.
  - o Binary.
  - o Nominal.
  - o Ordinal.
  - o Continuous.

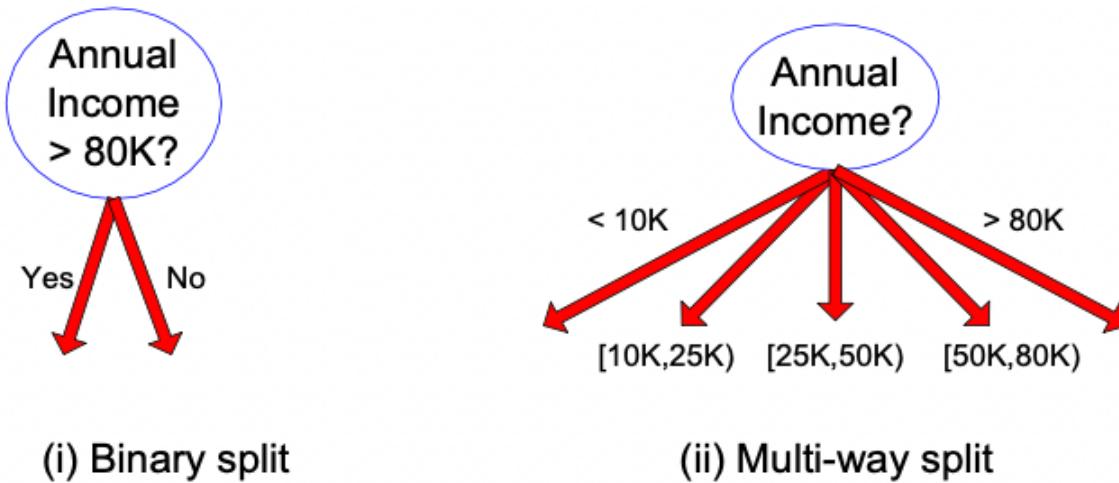
### Test Condition for Nominal Attributes

- **Multi-way split:**
  - o Use as many partitions as distinct values.
- **Binary split:**
  - o Divides values into 2 subsets.

### Test Condition for Ordinal Attributes

- **Multi-way split:**
  - o Use as many partitions as distinct values.
- **Binary split:**
  - o Divides values into 2 subsets.
  - o Preserve order property among attribute values.

### Test Condition for Continuous Attributes



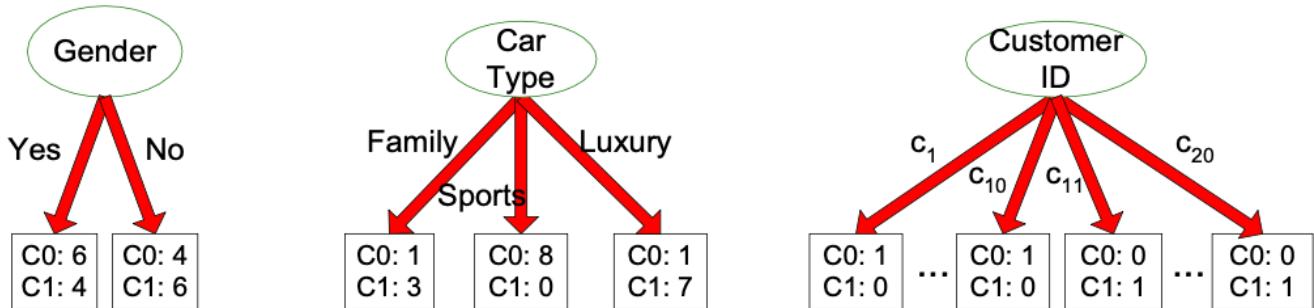
### Splitting Based on Continuous – Attributes

- Different ways of handling:
  - o **Discretization** to form an ordinal categorical attribute.
    - Ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.
      - Statistic: discretize once at beginning.
      - Dynamic: repeat at each node.
  - o **Binary Decision:**  $(A < v)$  or  $(A \geq v)$ .
    - Consider all possible splits and finds the best cut.
    - Can be more compute intensive.

## How to determine the Best Split

Customer Id	Gender	Car Type	Shirt Size	Class
1	M	Family	Small	C0
2	M	Sports	Medium	C0
3	M	Sports	Medium	C0
4	M	Sports	Large	C0
5	M	Sports	Extra Large	C0
6	M	Sports	Extra Large	C0
7	F	Sports	Small	C0
8	F	Sports	Small	C0
9	F	Sports	Medium	C0
10	F	Luxury	Large	C0
11	M	Family	Large	C1
12	M	Family	Extra Large	C1
13	M	Family	Medium	C1
14	M	Luxury	Extra Large	C1
15	F	Luxury	Small	C1
16	F	Luxury	Small	C1
17	F	Luxury	Medium	C1
18	F	Luxury	Medium	C1
19	F	Luxury	Medium	C1
20	F	Luxury	Large	C1

- Before Splitting: 10 records of class 0, 10 records of class 1.



- What test condition is the best?
- Greedy approach:
  - o Nodes with **purer** class distribution are preferred.
- Need a measure of node impurity.

C0: 5  
C1: 5

High degree of impurity

C0: 9  
C1: 1

Low degree of impurity

## Measures of Node Impurity

### I Gini Index

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

Where  $p_i(t)$  is the frequency of class  $i$  at node  $t$ , and  $c$  is the total number of classes

### I Entropy

$$Entropy = - \sum_{i=0}^{c-1} p_i(t) \log_2 p_i(t)$$

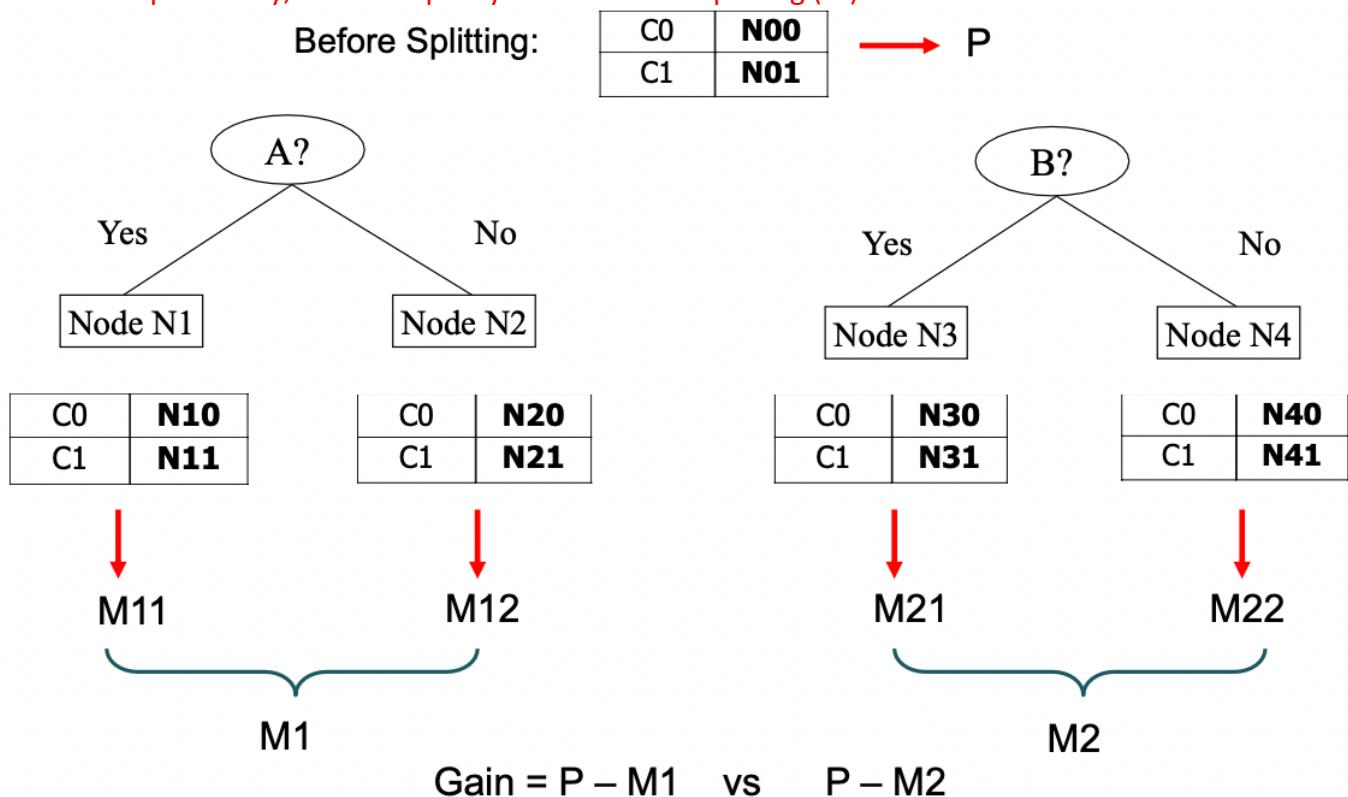
### I Misclassification error

$$Classification\ error = 1 - \max[p_i(t)]$$

## Finding the Best Split

1. Compute impurity measure ( $P$ ) before splitting.
2. Compute impurity measure ( $M$ ) after splitting.
  - o Compute impurity measure of each child node.
  - o  $M$  is the weighted impurity of child nodes.
3. Choose the attribute test condition that produces the highest gain.
  - o Gain =  $P - M$

Or equivalently, lowest impurity measure after splitting ( $M$ ).



## Measure of Impurity

- Gini Index for a given node  $t$

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

Where  $p_i(t)$  is the frequency of class  $i$  at node  $t$ , and  $c$  is the total number of classes

- Maximum of  $1-1/c$  when records are equally distributed among all classes, implying the least beneficial situation for classification.
- Minimum of 0 when all records belong to one class, implying the most beneficial situation for classification.
- Gini index is used in decision tree algorithms such as CART, SPRINT.
- For 2-class problem ( $p$ ,  $1-p$ ):
  - o GINI =  $1 - p^2 - (1-p)^2 = 2p(1-p)$

C1	<b>0</b>
C2	<b>6</b>
<b>Gini=0.000</b>	

C1	<b>1</b>
C2	<b>5</b>
<b>Gini=0.278</b>	

C1	<b>2</b>
C2	<b>4</b>
<b>Gini=0.444</b>	

C1	<b>3</b>
C2	<b>3</b>
<b>Gini=0.500</b>	

### Computing GINI Index of a Single Node

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

C1	<b>0</b>
C2	<b>6</b>

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$\text{Gini} = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$$

C1	<b>1</b>
C2	<b>5</b>

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$\text{Gini} = 1 - (1/6)^2 - (5/6)^2 = 0.278$$

C1	<b>2</b>
C2	<b>4</b>

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$\text{Gini} = 1 - (2/6)^2 - (4/6)^2 = 0.444$$

## Computing GINI Index for a Collection of Nodes

When a node  $p$  is split into  $k$  partitions (children)

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

where,  $n_i$  = number of records at child  $i$ ,

$n$  = number of records at parent node  $p$ .

### Binary Attributes: Computing GINI Index

- Splits into 2 partitions (child nodes).
- Effect of Weight partitions:
  - o Larger and purer partitions are sought.
- $Gini(N1) = 1 - (5/6)^2 - (1/6)^2 = 0.278$
- $Gini(N2) = 1 - (2/6)^2 - (4/6)^2 = 0.444$

$$\begin{aligned} & \text{Weighted Gini of N1 N2} \\ & = 6/12 * 0.278 + \\ & \quad 6/12 * 0.444 \\ & = 0.361 \end{aligned}$$

$$\text{Gain} = 0.486 - 0.361 = 0.125$$

```

graph TD
    B((B?)) -- Yes --> N1[Node N1]
    B -- No --> N2[Node N2]
  
```

	Parent
C1	7
C2	5
<b>Gini</b>	<b>0.486</b>

	N1	N2
C1	5	2
C2	1	4
<b>Gini</b>	<b>0.361</b>	

### Categorical Attributes: Computing Gini Index

- For each distinct value, gather counts for each class in the dataset.
- Use the count matrix to make decisions.

Multi-way split

CarType			
	Family	Sports	Luxury
C1	1	8	1
C2	3	0	7
<b>Gini</b>	<b>0.163</b>		

Two-way split  
(find best partition of values)

CarType		
	{Sports, Luxury}	{Family}
C1	9	1
C2	7	3
<b>Gini</b>	<b>0.468</b>	

CarType		
	{Sports}	{Family, Luxury}
C1	8	2
C2	0	10
<b>Gini</b>	<b>0.167</b>	

- Which of these is the best?

## Continuous Attributes: Computing Gini Index

- Use Binary Decision based on one value.
- Several Choices for the splitting value.
  - o Number of possible splitting values = Number of distinct values.
- Each Splitting value has a count matrix associate with it.
  - o Class counts in each of the partitions,  $A \leq v$  and  $A > v$ .
- Simple method to choose the best v:
  - o For each v, scan the database to gather count matrix and compute its Gini index.
  - o Computationally inefficient! Repetition of work.
- For efficient computation: For each attribute,
  - o Sort the attribute on values.
  - o Linearly scan these values, each time updating the count matrix and computing GINI index.
  - o Choose the split position that has the least GINI index.

ID	Home Owner	Marital Status	Annual Income	Defaulted
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

Cheat	No	No	No	Yes	Yes	Yes	No	No	No	No	No
Annual Income											
Sorted Values	→	60	70	75	85	90	95	100	120	125	220
Split Positions	→	55	65	72	80	87	92	97	110	122	172
	<=	>	<=	>	<=	>	<=	>	<=	>	<=
Yes	0	3	0	3	0	3	1	2	2	1	3
No	0	7	1	6	2	5	3	4	3	4	4
Gini	0.420	0.400	0.375	0.343	0.417	0.400	0.300	0.343	0.375	0.400	0.420

## Measure of Impurity: Entropy

- At each node t:

$$Entropy = - \sum_{i=0}^{c-1} p_i(t) \log_2 p_i(t)$$

- Where  $p_i(t)$  is the frequency of class i at node t, and c is total number of classes.
  - o Maximum of  $\log_2 c$  when records are equally distributed among all classes, implying the least beneficial situation for classification.
  - o Minimum of 0 when all records belong to one class, implying most beneficial situation of classification.
- Entropy based computations are quite like the GINI index computations.

## Computing Entropy of a Single Node

$$Entropy = - \sum_{i=0}^{c-1} p_i(t) \log_2 p_i(t)$$

C1	<b>0</b>
C2	<b>6</b>

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$\text{Entropy} = -0 \log 0 - 1 \log 1 = -0 - 0 = 0$$

C1	<b>1</b>
C2	<b>5</b>

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$\text{Entropy} = -(1/6) \log_2 (1/6) - (5/6) \log_2 (1/6) = 0.65$$

C1	<b>2</b>
C2	<b>4</b>

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$\text{Entropy} = -(2/6) \log_2 (2/6) - (4/6) \log_2 (4/6) = 0.92$$

## Computing Information Gain After Splitting

- Information Gain:

$$Gain_{split} = Entropy(p) - \sum_{i=1}^k \frac{n_i}{n} Entropy(i)$$

- o Parent Node, p is split into k partitions (children).
- o  $n_i$  is number of records in child node i.
- Choose the split that achieve most reduction (maximizes GAIN).
- Used in the ID3 and C4.5 decision tree algorithm.
- Information gain is the mutual information between the class variable and the splitting variable.

## Problem with large number of partitions

- Node impurity measures tend to prefer splits that result in large number of partitions, each being small but pure.
- Customer ID has highest information gain because entropy for all the children is zero.

## Gain Ratio

$$Gain\ Ratio = \frac{Gain_{split}}{Split\ Info} \quad Split\ Info = - \sum_{i=1}^k \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

- o Parent Node, p is split into k partitions (children).
- o  $n_i$  is the number of records in child node i.
- Adjusts Information Gain by the entropy of the partitioning (Split Info).

- Higher entropy partitioning (large number of small partitions) is penalized.
- Used in C4.5 algorithm.
- Designed to overcome the disadvantage of Information Gain.

CarType		
	Family	Sports
C1	1	8
C2	3	0
Gini	<b>0.163</b>	

$$\text{SplitINFO} = 1.52$$

CarType		
	{Sports, Luxury}	{Family}
C1	9	1
C2	7	3
Gini	<b>0.468</b>	

$$\text{SplitINFO} = 0.72$$

CarType		
	{Sports}	{Family, Luxury}
C1	8	2
C2	0	10
Gini	<b>0.167</b>	

$$\text{SplitINFO} = 0.97$$

### Measure of Impurity: Classification Error

- Classification error at a node t:

$$Error(t) = 1 - \max_i [p_i(t)]$$

- Maximum of  $1 - 1/c$  when records are equally distributed among all classes, implying the least interesting situation.
- Minimum of 0 when all records belong to one class, implying the most interesting situation.

### Computing Error of a Single Node

$$Error(t) = 1 - \max_i [p_i(t)]$$

C1	<b>0</b>
C2	<b>6</b>

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$\text{Error} = 1 - \max(0, 1) = 1 - 1 = 0$$

C1	<b>1</b>
C2	<b>5</b>

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$\text{Error} = 1 - \max(1/6, 5/6) = 1 - 5/6 = 1/6$$

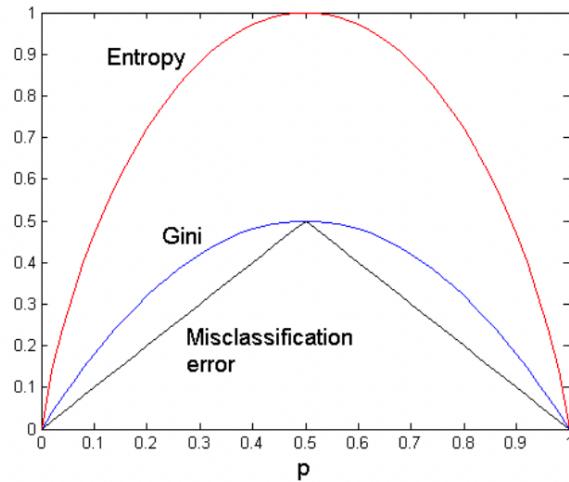
C1	<b>2</b>
C2	<b>4</b>

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

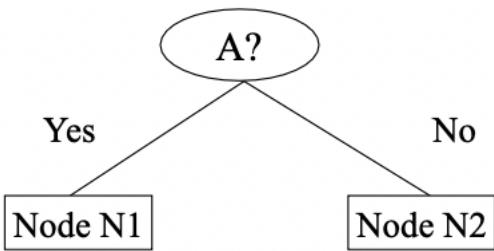
$$\text{Error} = 1 - \max(2/6, 4/6) = 1 - 4/6 = 1/3$$

### Comparison among Impurity Measures

For a 2-class problem:



## Misclassification Error vs. Gini Index



	<b>Parent</b>
C1	<b>7</b>
C2	<b>3</b>
<b>Gini = 0.42</b>	

	<b>N1</b>	<b>N2</b>
C1	<b>3</b>	<b>4</b>
C2	<b>0</b>	<b>3</b>
<b>Gini=0.342</b>		

- Gini improves but error maintains the same.

	<b>N1</b>	<b>N2</b>
C1	<b>3</b>	<b>4</b>
C2	<b>0</b>	<b>3</b>
<b>Gini=0.342</b>		

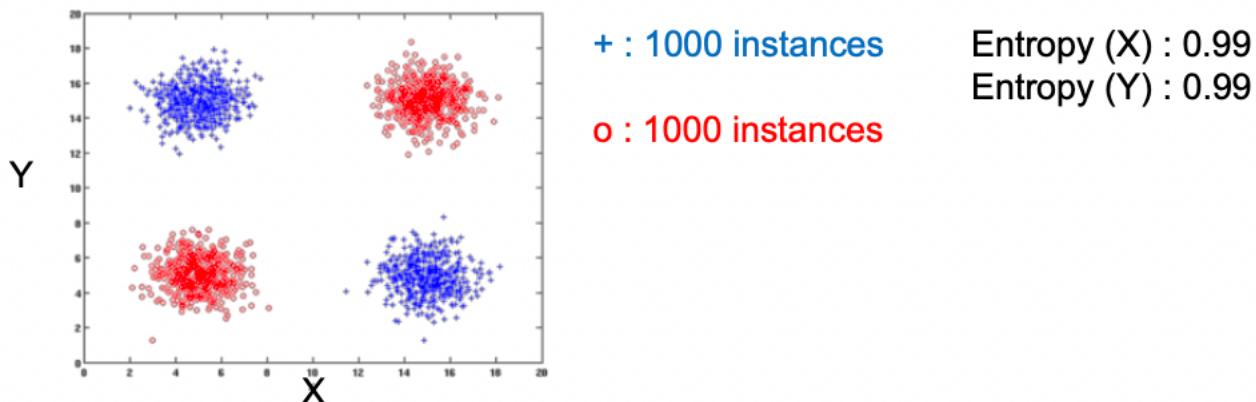
	<b>N1</b>	<b>N2</b>
C1	<b>3</b>	<b>4</b>
C2	<b>1</b>	<b>2</b>
<b>Gini=0.416</b>		

- Misclassification error for all three cases = 0.3!

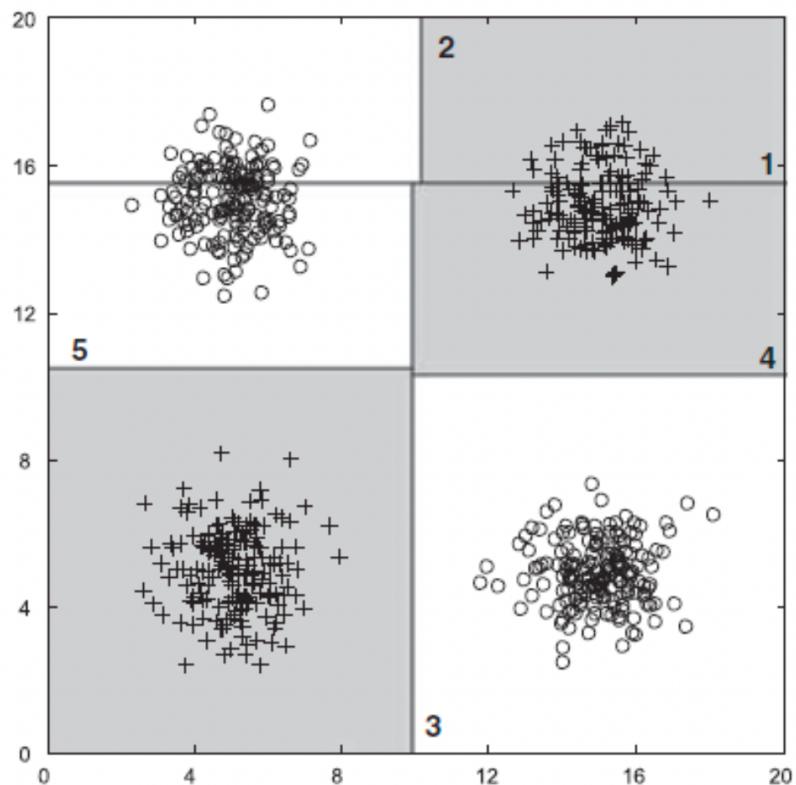
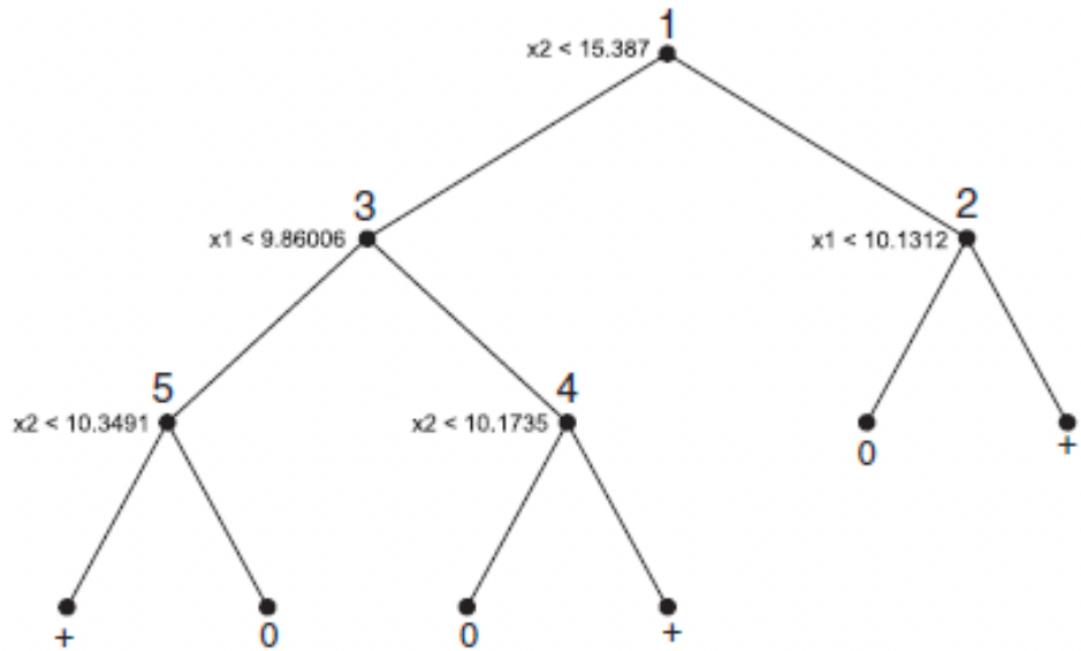
## Decision Tree Based Classification

- Advantages:
  - o Relatively inexpensive to construct.
  - o Extremely fast at classifying unknown records.
  - o Easy to interpret for small-sized trees.
  - o Robust to noise (especially when methods to avoid overfitting are employed).
  - o Can easily handle redundant attributes.
  - o Can easily handle irrelevant attributes (unless the attributes are interacting).
- Disadvantages:
  - o Due to the greedy nature of splitting criterion, interacting attributes (that can distinguish between classes together but not individually) may be passed over in favor of other attributes that are less discriminating.
  - o Each decision boundary involves only a single attribute.

## Handling Interactions

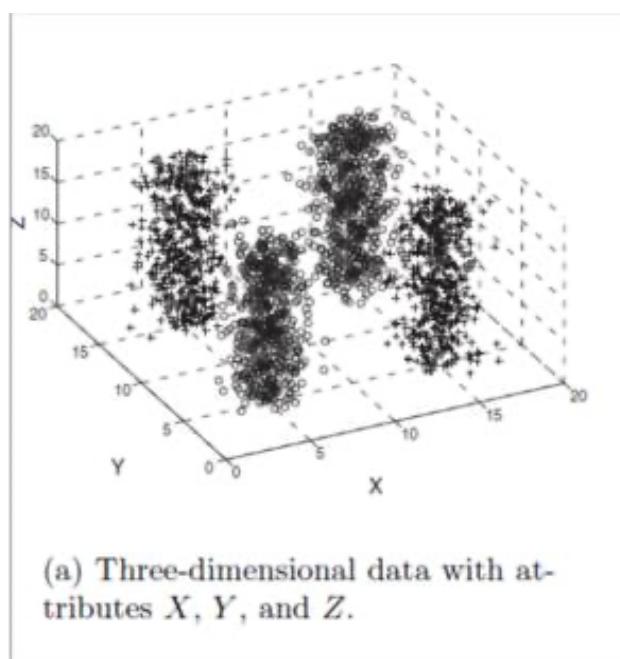


## Handling Interactions – 6 leaf nodes



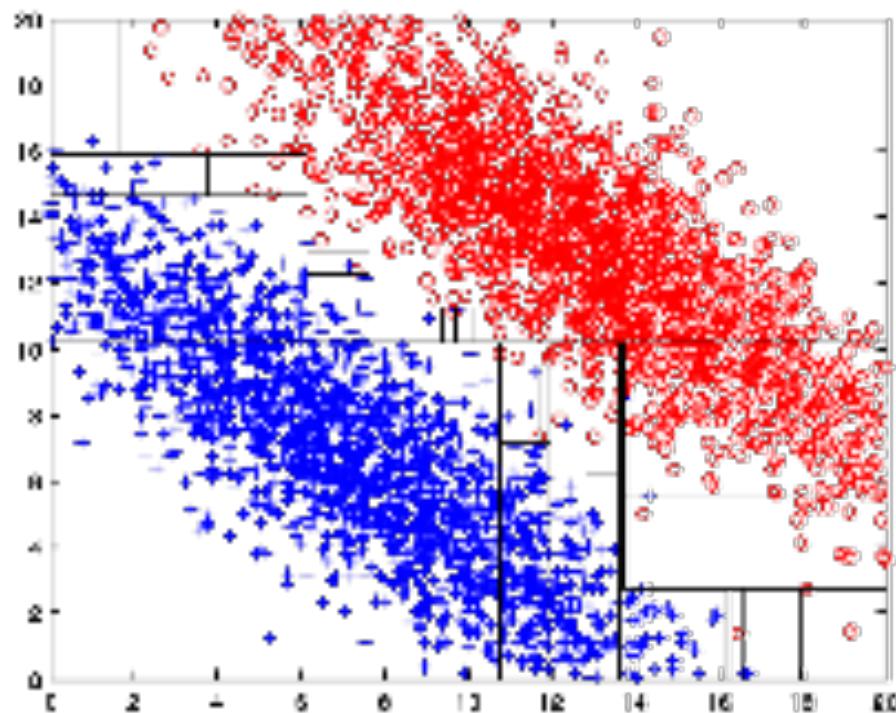
## Handling Interactions Given Irrelevant Attributes

- Adding Z as a noisy attribute generated from a uniform distribution.
- Attribute Z will be chosen for splitting.



Entropy ( $X$ ) : 0.99  
Entropy ( $Y$ ) : 0.99  
Entropy ( $Z$ ) : 0.98

- Both positive (+) and negative (o) classes generated from skewed Gaussians with centers at (8, 8) and (12, 12) respectively.



## Table of Contents

<i>Classification: Definition</i> .....	1
<i>General Approach for Building Classification Model</i> .....	1
<i>Classification Techniques</i> .....	2
<i>Example of a Decision Tree</i> .....	2
<i>Decision Tree Classification Task</i> .....	2
<i>Decision Tree Induction</i> .....	3
<i>General Structure of Hunt's Algorithm</i> .....	3
<i>Hunt's Algorithm</i> .....	3
<i>Design Issues of Decision Tree Induction</i> .....	3
<i>Methods for Expressing Test Conditions</i> .....	4
<i>Test Condition for Nominal Attributes</i> .....	4
<i>Test Condition for Ordinal Attributes</i> .....	4
<i>Test Condition for Continuous Attributes</i> .....	4
<i>Splitting Based on Continuous – Attributes</i> .....	4
<i>How to determine the Best Split</i> .....	5
<i>Measures of Node Impurity</i> .....	6
<i>Finding the Best Split</i> .....	6
<i>Measure of Impurity</i> .....	7
<i>Computing GINI Index of a Single Node</i> .....	7
<i>Computing GINI Index for a Collection of Nodes</i> .....	8
<i>Binary Attributes: Computing GINI Index</i> .....	8
<i>Categorical Attributes: Computing Gini Index</i> .....	8
<i>Continuous Attributes: Computing Gini Index</i> .....	9
<i>Measure of Impurity: Entropy</i> .....	9
<i>Computing Entropy of a Single Node</i> .....	10
<i>Computing Information Gain After Splitting</i> .....	10
<i>Problem with large number of partitions</i> .....	10
<i>Gain Ratio</i> .....	10
<i>Measure of Impurity: Classification Error</i> .....	11
<i>Computing Error of a Single Node</i> .....	11
<i>Comparison among Impurity Measures</i> .....	11
<i>Misclassification Error vs. Gini Index</i> .....	12
<i>Decision Tree Based Classification</i> .....	12
<i>Handling Interactions</i> .....	12
<i>Handling Interactions – 6 leaf nodes</i> .....	13
<i>Handling Interactions Given Irrelevant Attributes</i> .....	14

# Data Mining: Cluster Analysis: Basic Concepts and Methods

## What is Cluster Analysis?

- Cluster: A collection of data objects.
  - o Similar (or related) to one another within the same group.
  - o Dissimilar (or unrelated) to the objects in other groups.
- Cluster Analysis (or clustering, data segmentation, etc.):
  - o Finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters.
- **Unsupervised learning:** no predefined classes (i.e., *learning by observations* vs. *learning by examples*: supervised).
- Typical applications:
  - o As a **stand-alone tool** to get insight into data distribution.
  - o As a **preprocessing** step for other algorithms.

## Applications of Cluster Analysis

- Data reduction:
  - o Summarization: Preprocessing for regression, PCA, classification, and association analysis.
  - o Compression: Image processing: vector quantization.
- Hypothesis generation and testing.
- Prediction based on groups.
  - o Cluster & find characteristics/patterns for each group.
- Finding K-nearest Neighbors.
  - o Localizing search to one or a small number of clusters.
- Outlier detection: Outliers are often viewed as those “far away” from any cluster.

## Clustering: Application Examples

- Biology: taxonomy of living things: kingdom, phylum, class, order, family, genus, and species.
- Information retrieval: document clustering.
- Land use: Identification of areas of similar land use in an earth observation database.
- Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs.
- City-planning: Identifying groups of houses according to their house type, value, and geographical location.
- Earthquake studies: Observed earthquake epicenters should be clustered along continent faults.
- Climate: understanding earth climate, find patterns of atmospheric and ocean.
- Economic Science: market research.

## Basic Steps to Develop a Cluster Task

- Feature Selection:
  - o Select info concerning the task of interest.
  - o Minimal information redundancy.
- Proximity measure:
  - o Similarity of two feature vectors.
- Clustering criterion:
  - o Expressed via a cost function or some rules.
- Clustering algorithms:
  - o Choice of algorithms.
- Validation of the results:
  - o Validation test (also, **clustering tendency** test).
- Interpretation of the results:
  - o Integration with applications.

## Quality: What is Good Clustering?

- A good clustering method will produce high quality clusters:
  - o High intra-class similarity: **cohesive** within clusters.
  - o Low inter-class similarity: **distinctive** between clusters.
- The quality of a clustering method depends on:
  - o The similarity measure used by the method.
  - o Its implementation, and
  - o Its ability to discover some or all the hidden patterns.

## Measure the Quality of Clustering

- **Dissimilarity/Similarity metric:**
  - o Similarity is expressed in terms of a distance function, typically metric:  $d(i, j)$ .
  - o The definitions of **distance functions** are usually rather different for interval-scaled, Boolean, categorical, ordinal ratio, and vector variables.
  - o Weights should be associated with different variables based on applications and data semantics.
- Quality of clustering:
  - o There is usually a separate “quality” function that measures the “goodness” of a cluster.
  - o It is hard to define “similar enough” or “good enough”. The answer is typically highly subjective.

## Considerations for Cluster Analysis

- Partitioning criteria:
  - o Single level.
  - o Hierarchical partitioning (often, multi-level hierarchical partitioning is desirable).
- Separation of clusters:
  - o Exclusive (e.g., one customer belongs to only one region).
  - o Non-exclusive (e.g., one document may belong to more than one class).
- Similarity measure:
  - o Distance-based (Euclidian, road network, vector).
  - o Connectivity-based (density, contiguity).
- Clustering space:
  - o Full space (often when low dimensional).
  - o Subspaces (often in high-dimensional clustering).

## Requirements and Challenges

- Scalability:
  - o Clustering all the data instead of only on samples.
- Ability to deal with different types of attributes:
  - o Numerical, binary, categorical, ordinal, linked, and mixture of these.
- Constraint-based clustering:
  - o User may give inputs on constraints.
  - o Use domain knowledge to determine input parameters.
- Interpretability and usability.
- Others:
  - o Discovery of clusters with arbitrary shape.
  - o Ability to deal with noisy data.
  - o Incremental clustering and insensitivity to input order.
  - o High dimensionality.

## Major Clustering Approaches

- Partitioning approach:
  - o Construct various partitions and then evaluate them by some criterion.
  - o E.g., minimizing the sum of square errors.
  - o Typical methods: k-means, k-medoids, CLARANS.
- Hierarchical approach:
  - o Create a hierarchical decomposition of the set of data (or objects) using some criterion.
  - o Typical methods: Diana, Agnes, BIRCH, CAMELEON.
- Density-based approach:
  - o Based on connectivity and density functions.
  - o Typical methods: DBSCAN, OPTICS, Den Clue.
- Grid-based approach:
  - o Based on a multiple-level granularity structure.
  - o Typical methods: STING, Wave Cluster, CLIQUE.
- Model-based:
  - o A model is hypothesized for each of the clusters and tries to find the best fit of that model to each other.
  - o Typical methods: EM, SOM, COBWEB.
- Frequent pattern-based:
  - o Based on the analysis of frequent patterns.
  - o Typical methods: p-Cluster.
- User-guided or constraint-based:
  - o Clustering by considering user-specified or application-specific constraints.
  - o Typical methods: COD (obstacles), constrained clustering.
- Link-based clustering:
  - o Objects are often linked together in various ways.
  - o Massive links can be used to cluster objects: SimRank, LinkClus.

## Partitioning Algorithms: Basic Concept

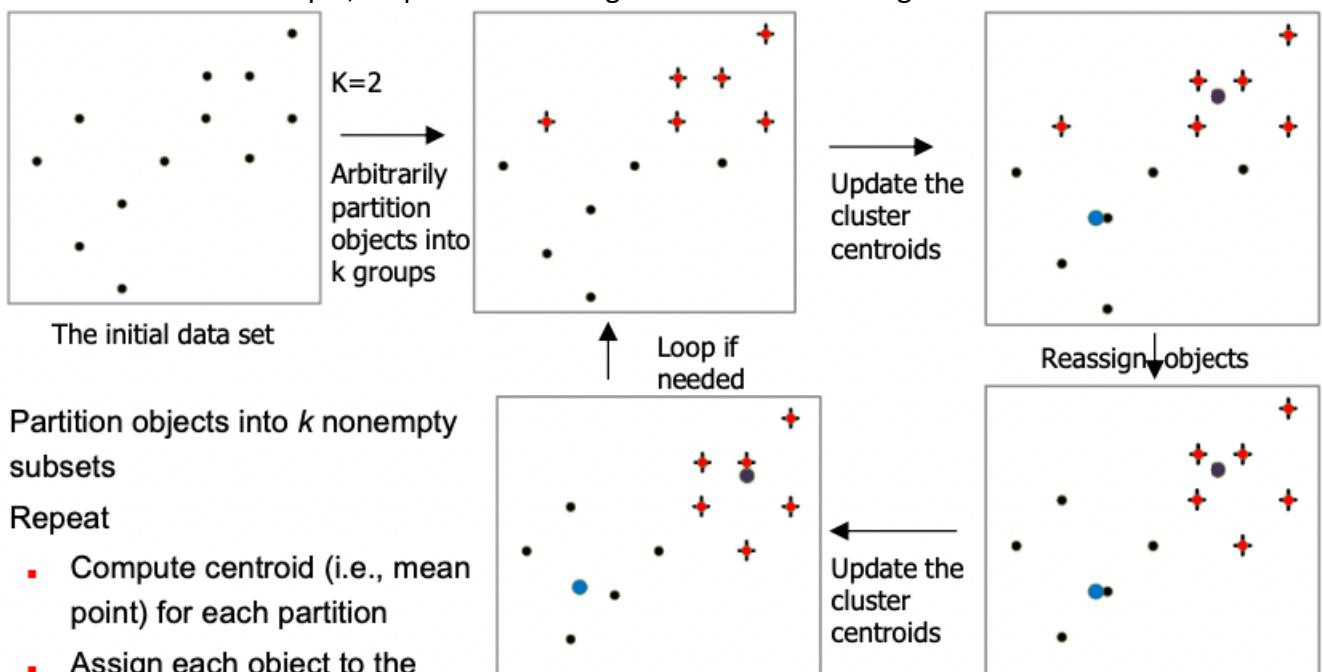
- **Partitioning method:** Partitioning a database D of n objects into a set of k clusters, such that the sum of squared distances is minimized (where  $c_i$  is the centroid of medoid of cluster  $C_i$ ).

$$E = \sum_{i=1}^k \sum_{p \in C_i} (d(p, c_i))^2$$

- Given k, find a partition of k clusters that optimized the chosen partitioning criterion.
  - o Global optimal: exhaustively enumerate all partitions.
  - o Heuristic methods: k-means and k-medoids algorithms.
  - o K-means (MacQueen'67, Klogd'57/'82): Each cluster is represented by the center of the cluster.
  - o K-medoids or PAM (Partition Around Medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster.

## The K-Means Clustering Method.

- Given k, the k-means algorithm is implemented in four steps:
  - o Partition objects into k non-empty subsets.
  - o Compute seed points as the centroids of the clusters of the current partitioning (the centroid is the center, i.e., **mean point** of the cluster).
  - o Assign each object to the cluster with the nearest seed point.
  - o Go back to step 2, stop when the assignment does not change.

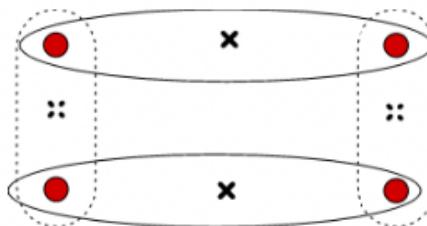


## Comments on the K-means Method

- Strength: Efficient  $O(t*k*n)$ , where  $n$  is number of objects,  $k$  is number of clusters, and  $t$  is number of iterations. Normally  $k, t \ll n$ .
  - o PAM:  $O(k(n-k)^2)$
  - o CLARA:  $O(ks^2 + k(n-k))$
- Comment: Often terminates at a local optimal.
- Weakness:
  - o Applicable only to objects in a continuous  $n$ -dimensional space.
    - Using the k-modes method for categorical data.
    - In comparison, k-medoids can be applied to a wide range of data.
  - o Need to specify  $k$ : the number of clusters, in advance (there are ways to automatically determine the best  $k$  (see Hastie et al., 2009)).
  - o Sensitive to noisy data and outliers.
  - o Not suitable to discover clusters with non-convex shapes.

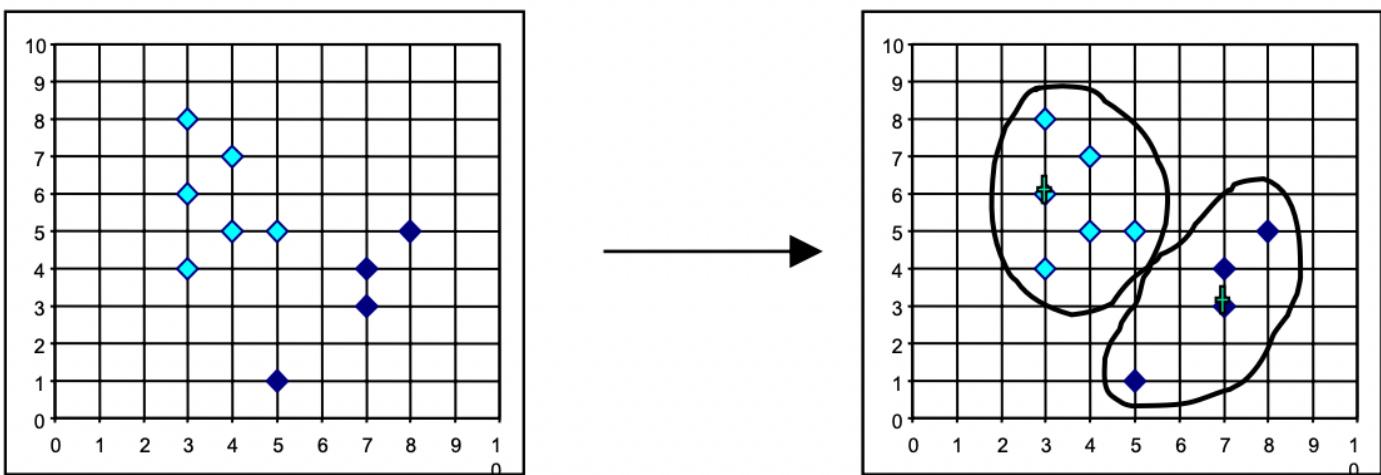
## Variations of the K-Means Method

- Most of the variants of the k-means which differ in:
  - o Selection of the initial  $k$  means.
  - o Dissimilarity calculations.
  - o Strategies to calculate cluster means.
- Handling categorical data: k-modes.
  - o Replacing means of clusters with modes.
  - o Using new dissimilarity measures to deal with categorical objects.
  - o Using a frequency-based method to update modes of clusters.
  - o A mixture of categorical and numerical data: k-prototype method.

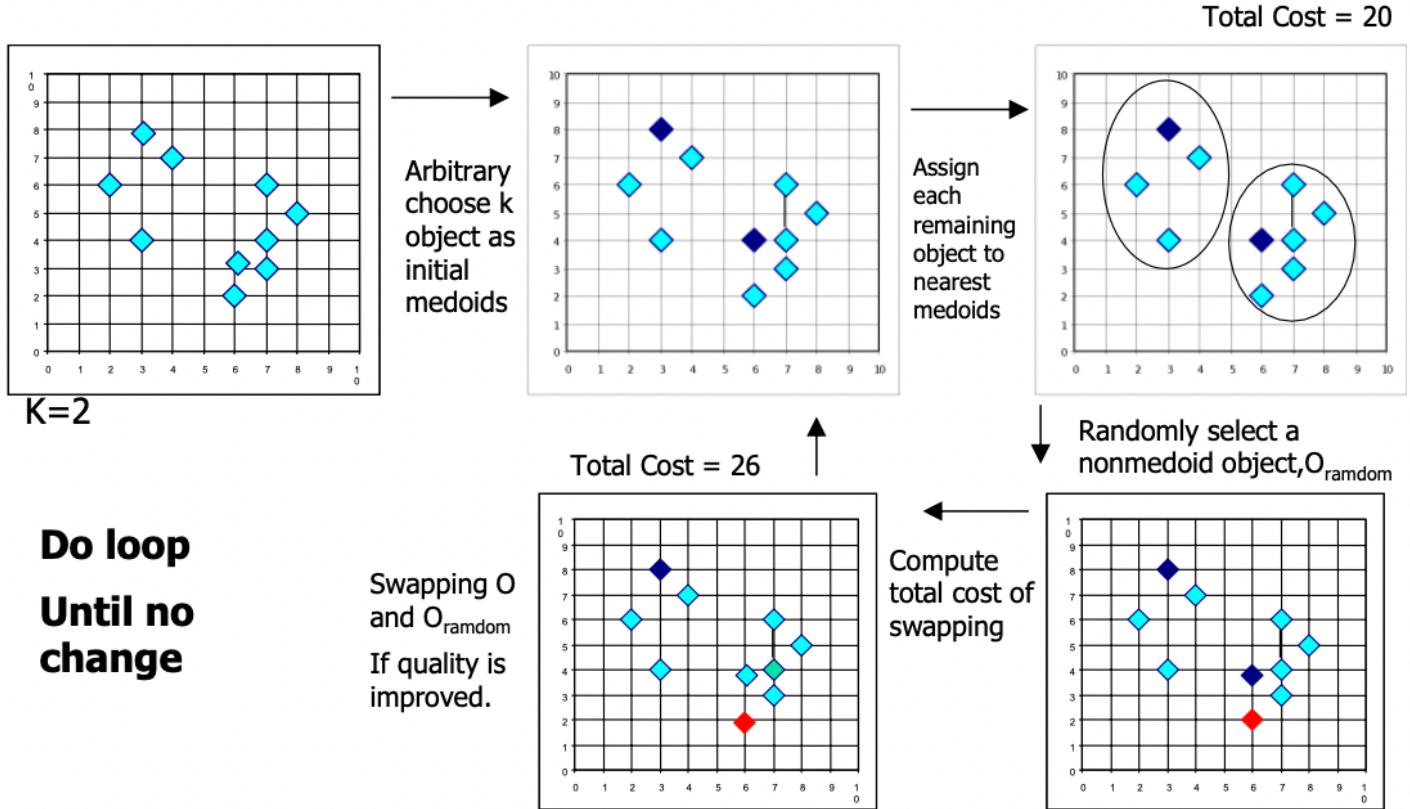


## What is the Problem of the K-Means Method

- The k-means algorithm is sensitive to outliers.
  - o Since an object with an extremely large value may substantially distort the distribution of the data.
- K-Medoids: Instead of taking the **mean** value of the object in a cluster as a reference point, **medoids** can be used, which is the most **centrally located** object in a cluster.

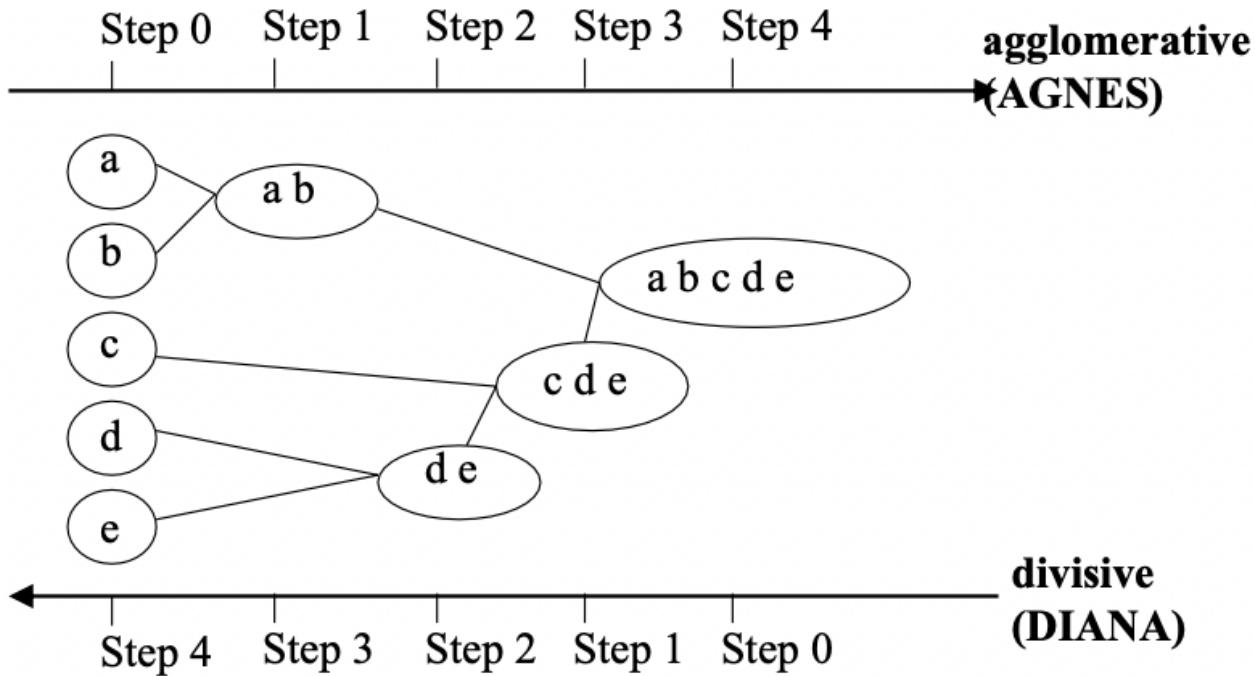


## PAM: A Typical K-Medoids Algorithm



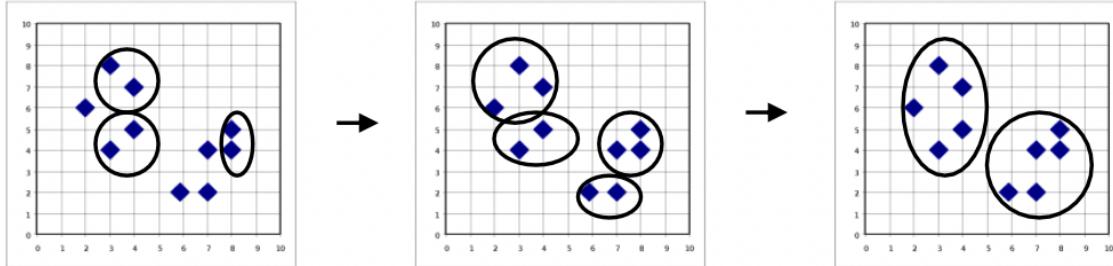
## Hierarchical Clustering

- Use distance matrix as clustering criteria. This method does not require the number of clusters k as an input but needs a termination condition.

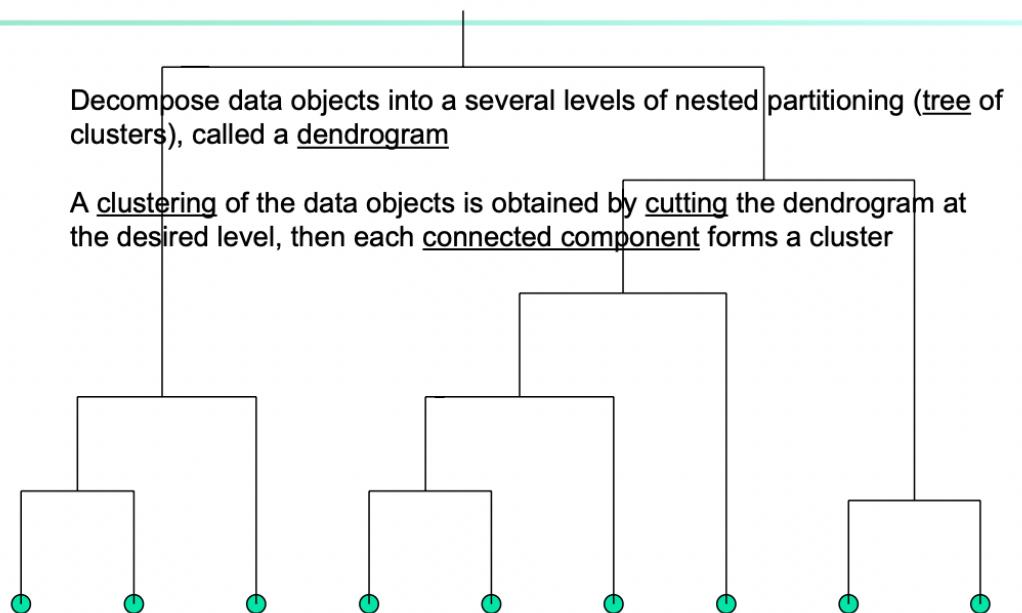


## AGNES (Agglomerative Nesting)

- Introduced in Kaufmann and Rousseau (1990).
- Implemented in statistical packages, e.g., Splus.
- Use the single-link method and the dissimilarity matrix.
- **Merge nodes that have the least dissimilarity.**
- Go on in a non-descending fashion.
- Eventually all nodes belong to the same cluster.

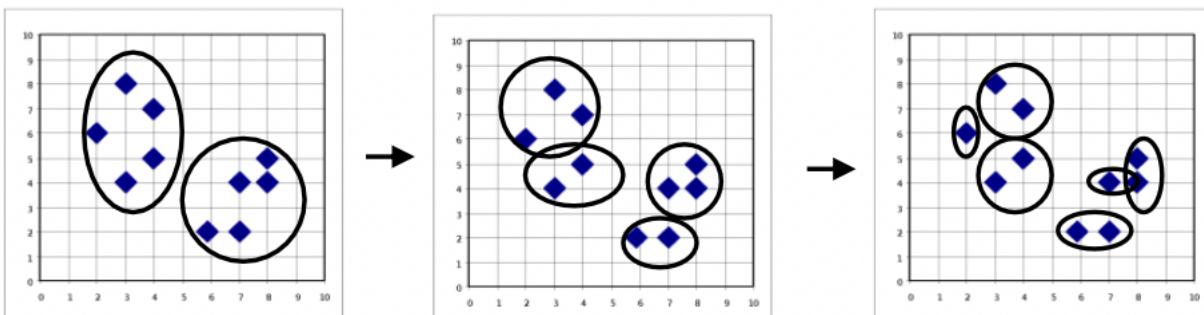


Dendrogram: Shows how Clusters Are Merged



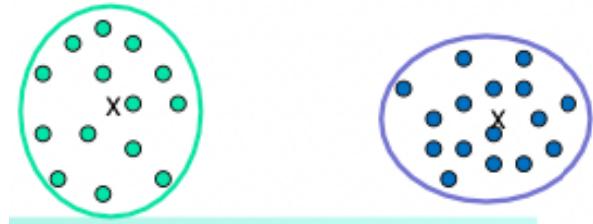
## DIANA (Divisive Analysis)

- Introduced in Kaufmann and Rousseau (1990).
- Implemented in statistical analysis package, e.g., Splus.
- Inverse order of AGNES.
- Eventually each node forms a cluster on its own.



## Distance Between Clusters

- Single link: smallest distance between an element in one cluster and an element in the other:
  - o  $\text{dist}(K_i, K_j) = \min(t_{ip}, t_{jq})$
- Complete link: largest distance between an element in one cluster and an element in the other:
  - o  $\text{dist}(K_i, K_j) = \max(t_{ip}, t_{jq})$
- Average: avg distance between an element in one cluster and an element in the other:
  - o  $\text{dist}(K_i, K_j) = \text{avg}(t_{ip}, t_{jq})$
- Centroid: distance between the centroids of 2 clusters:
  - o  $\text{dist}(K_i, K_j) = \text{dist}(C_i, C_j)$
- Medoid: distance between the medoids of two clusters:
  - o  $\text{dist}(K_i, K_j) = \text{dist}(M_i, M_j)$
  - o Medoid: a chosen, centrally located object in the cluster.



## Centroid, Radius, and Diameter of a Cluster (for numerical data sets)

- Centroid: the “middle” of a cluster.
- $$C_m = \frac{\sum_{i=1}^N (t_{ip})}{N}$$
- Radius: Square root of average distance from any point of the cluster to its centroid.
- $$R_m = \sqrt{\frac{\sum_{i=1}^N (t_{ip} - c_m)^2}{N}}$$
- Diameter: Square root of average mean squared distance between all pairs of points in the cluster.

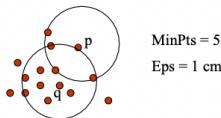
$$D_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (t_{ip} - t_{jq})^2}{N(N-1)}}$$

## Density-Based Clustering Methods

- Clustering based on density (local cluster criterion), such as density-connected points.
- Major features:
  - o Discover clusters of arbitrary shape.
  - o Handle noise.
  - o One scan.
  - o Need density parameters as termination condition.
- Several interesting studies:
  - **DBSCAN**: Ester, et al. (KDD’96)
  - **OPTICS**: Ankerst, et al (SIGMOD’99).
  - **DENCLUE**: Hinneburg & D. Keim (KDD’98)
  - **CLIQUE**: Agrawal, et al. (SIGMOD’98) (more grid-based)

## Density-Based Clustering: Basic Concepts

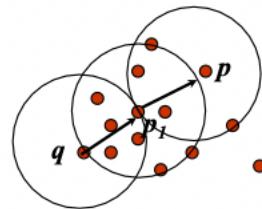
- 2 parameters:
  - o **Eps**: Maximum radius of the neighborhood.
  - o **MinPts**: Minimum number of points in an Eps-neighborhood of that point.
- $N_{Eps}(q) = \{p \text{ belongs to } D \mid \text{dist}(p, q) \leq Eps\}$
- **Directed density-reachable**: A point  $p$  is directly density-reachable from a point  $q$  w.r.t Eps, MinPts if:
  - o  $p$  belongs to  $N_{Eps}(q)$
  - o core point condition:
$$|N_{Eps}(q)| \geq \text{MinPts}$$



## Density-Reachable and Density-Connected

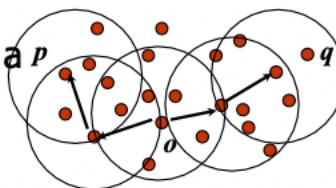
### Density-reachable:

- A point  $p$  is **density-reachable** from a point  $q$  w.r.t.  $Eps$ ,  $MinPts$  if there is a chain of points  $p_1, \dots, p_n$ ,  $p_1 = q$ ,  $p_n = p$  such that  $p_{i+1}$  is directly density-reachable from  $p_i$ .



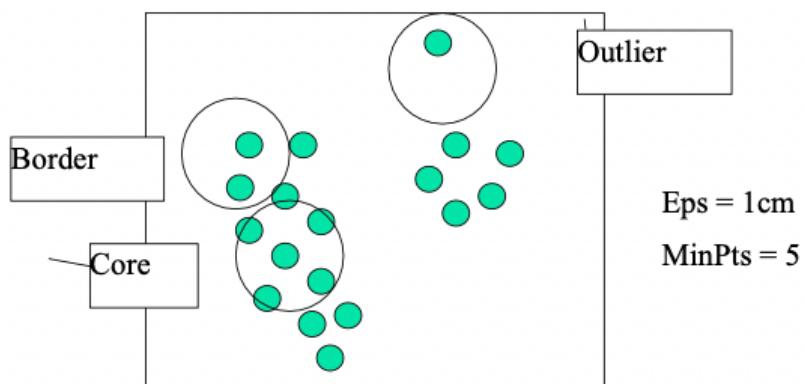
### Density-connected

- A point  $p$  is **density-connected** to a point  $q$  w.r.t.  $Eps$ ,  $MinPts$  if there is a point  $o$  such that both,  $p$  and  $q$  are density-reachable from  $o$  w.r.t.  $Eps$  and  $MinPts$



## DBSCAN: Density-Based Spatial Clustering of Applications with Noise

- Relies on a *density-based* notion of cluster: A *cluster* is defined as a maximal set of density-connected points
- Discovers clusters of arbitrary shape in spatial databases with noise



## DBSCAN: Sensitive to Parameters

Figure 8. DBScan results for DS1 with MinPts at 4 and Eps at (a) 0.5 and (b) 0.4.

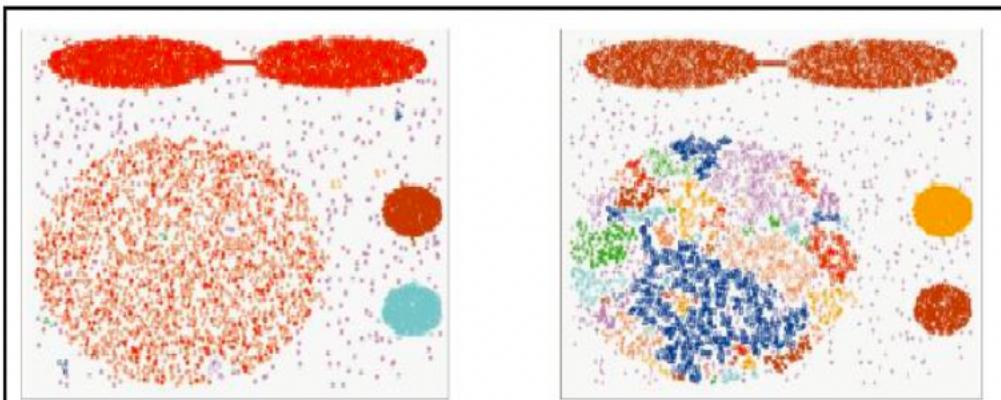
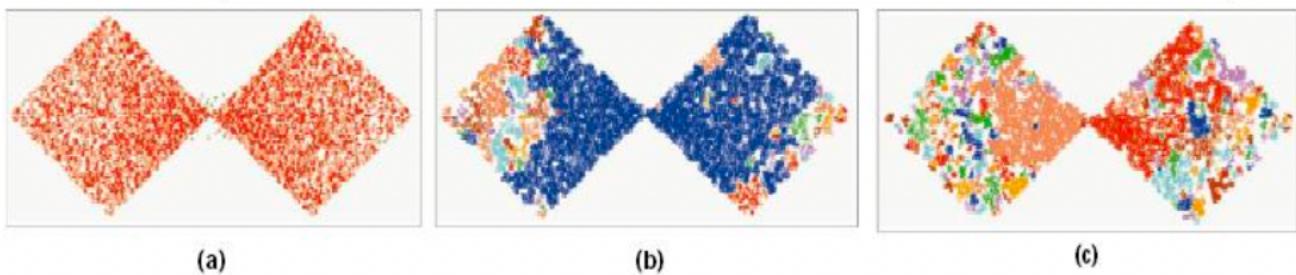


Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



## OPTICS: A Cluster-Ordering Method (1999)

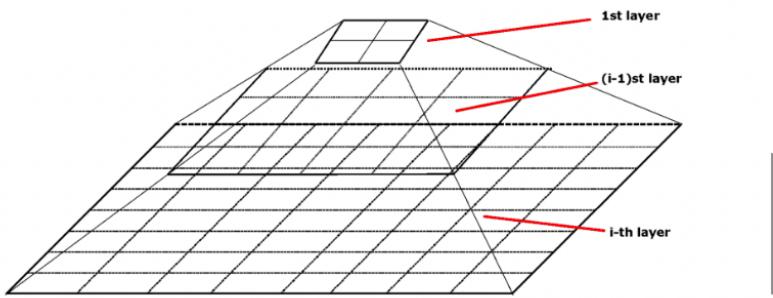
- **OPTICS: Ordering Points To Identify the Clustering Structure**
  - Ankerst, Breunig, Kriegel, and Sander (SIGMOD'99)
  - Produces a special order of the database wrt its density-based clustering structure
  - This cluster-ordering contains info equiv to the density-based clusterings corresponding to a broad range of parameter settings
  - Good for both automatic and interactive cluster analysis, including finding intrinsic clustering structure
  - Can be represented graphically or using visualization techniques

## Grid-Based Clustering Method

- Using multi-resolution grid data structure
- Several interesting methods
  - **STING** (a STatistical INformation Grid approach) by Wang, Yang and Muntz (1997)
  - **CLIQUE**: Agrawal, et al. (SIGMOD'98)
    - Both grid-based and subspace clustering
  - **WaveCluster** by Sheikholeslami, Chatterjee, and Zhang (VLDB'98)
    - A multi-resolution clustering approach using wavelet method

## STING: A Statistical Information Grid Approach

- Wang, Yang and Muntz (VLDB'97)
- The spatial area is divided into rectangular cells
- There are several levels of cells corresponding to different levels of resolution



## The STING Clustering Method

- Each cell at a high level is partitioned into a number of smaller cells in the next lower level
- Statistical info of each cell is calculated and stored beforehand and is used to answer queries
- Parameters of higher level cells can be easily calculated from parameters of lower level cell
  - *count, mean, s, min, max*
  - type of distribution—*normal, uniform*, etc.
- Use a top-down approach to answer spatial data queries
- Start from a pre-selected layer—typically with a small number of cells
- For each cell in the current level compute the confidence interval

## STING Algorithm and Its Analysis

- Remove the irrelevant cells from further consideration
- When finish examining the current layer, proceed to the next lower level
- Repeat this process until the bottom layer is reached
- Advantages:
  - Query-independent, easy to parallelize, incremental update
  - $O(K)$ , where  $K$  is the number of grid cells at the lowest level
- Disadvantages:
  - All the cluster boundaries are either horizontal or vertical, and no diagonal boundary is detected

## CLIQUE (Clustering In QUEst)

- Agrawal, Gehrke, Gunopulos, Raghavan (SIGMOD'98)
- Automatically identifying subspaces of a high dimensional data space that allow better clustering than original space
- CLIQUE can be considered as both density-based and grid-based
  - It partitions each dimension into the same number of equal length interval
  - It partitions an m-dimensional data space into non-overlapping rectangular units
  - A unit is dense if the fraction of total data points contained in the unit exceeds the input model parameter
  - A cluster is a maximal set of connected dense units within a subspace

## CLIQUE: The Major Steps

- Partition the data space and find the number of points that lie inside each cell of the partition.
- Identify the subspaces that contain clusters using the Apriori principle
- Identify clusters
  - Determine dense units in all subspaces of interests
  - Determine connected dense units in all subspaces of interests.
- Generate minimal description for the clusters
  - Determine maximal regions that cover a cluster of connected dense units for each cluster
  - Determination of minimal cover for each cluster

## Determine The Number Of Clusters

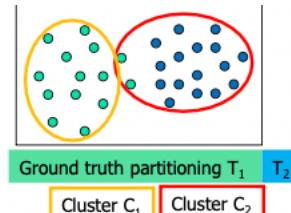
- Empirical method
  - # of clusters:  $k \approx \sqrt{n}/2$  for a dataset of  $n$  points, e.g.,  $n = 200$ ,  $k = 10$
- Elbow method
  - Use the turning point in the curve of sum of within cluster variance w.r.t the # of clusters
- Cross validation method
  - Divide a given data set into  $m$  parts
  - Use  $m - 1$  parts to obtain a clustering model
  - Use the remaining part to test the quality of the clustering
    - E.g., For each point in the test set, find the closest centroid, and use the sum of squared distance between all points in the test set and the closest centroids to measure how well the model fits the test set
  - For any  $k > 0$ , repeat it  $m$  times, compare the overall quality measure w.r.t. different  $k$ 's, and find # of clusters that fits the data the best

## Measuring Clustering Quality

- 3 kinds of measures: External, internal and relative
- External: supervised, employ criteria not inherent to the dataset
  - Compare a clustering against prior or expert-specified knowledge (i.e., the ground truth) using certain clustering quality measure
- Internal: unsupervised, criteria derived from data itself
  - Evaluate the goodness of a clustering by considering how well the clusters are separated, and how compact the clusters are, e.g., Silhouette coefficient
- Relative: directly compare different clusterings, usually those obtained via different parameter settings for the same algorithm

## Some Commonly Used External Measures

- Matching-based measures
  - Purity, maximum matching, F-measure
- Entropy-Based Measures
  - Conditional entropy, normalized mutual information (NMI), variation of information
- Pair-wise measures
  - Four possibilities: True positive (TP), FN, FP, TN
  - Jaccard coefficient, Rand statistic, Fowlkes-Mallow measure
- Correlation measures
  - Discretized Huber static, normalized discretized Huber static



## Table of Contents

<i>What is Cluster Analysis?</i> .....	1
<i>Applications of Cluster Analysis</i> .....	1
<i>Clustering: Application Examples</i> .....	1
<i>Basic Steps to Develop a Cluster Task</i> .....	2
<i>Quality: What is Good Clustering?</i> .....	2
<i>Measure the Quality of Clustering</i> .....	2
<i>Considerations for Cluster Analysis</i> .....	2
<i>Requirements and Challenges</i> .....	3
<i>Major Clustering Approaches</i> .....	3
<i>Partitioning Algorithms: Basic Concept</i> .....	4
<i>The K-Means Clustering Method</i> .....	4
<i>Comments on the K-means Method</i> .....	5
<i>Variations of the K-Means Method</i> .....	5
<i>What is the Problem of the K-Means Method</i> .....	5
<i>PAM: A Typical K-Medoids Algorithm</i> .....	6
<i>Hierarchical Clustering</i> .....	6
<i>AGNES (Agglomerative Nesting)</i> .....	7
<i>Dendrogram: Shows how Clusters Are Merged</i> .....	7
<i>DIANA (Divisive Analysis)</i> .....	7
<i>Distance Between Clusters</i> .....	8
<i>Centroid, Radius, and Diameter of a Cluster (for numerical data sets)</i> .....	8
<i>Density-Based Clustering Methods</i> .....	8
<i>Density-Based Clustering: Basic Concepts</i> .....	9
<i>Density-Reachable and Density-Connected</i> .....	9
<i>DBSCAN: Density-Based Spatial Clustering of Applications with Noise</i> .....	9
<i>DBSCAN: Sensitive to Parameters</i> .....	10
<i>OPTICS: A Cluster-Ordering Method (1999)</i> .....	10
<i>Grid-Based Clustering Method</i> .....	11
<i>STING: A Statistical Information Grid Approach</i> .....	11
<i>The STING Clustering Method</i> .....	11
<i>STING Algorithm and Its Analysis</i> .....	12
<i>CLIQUE (Clustering In QUEst)</i> .....	12
<i>CLIQUE: The Major Steps</i> .....	12

<i>Determine The Number Of Clusters</i> .....	<b>13</b>
<i>Measuring Clustering Quality</i> .....	<b>13</b>
<i>Some Commonly Used External Measures</i> .....	<b>13</b>

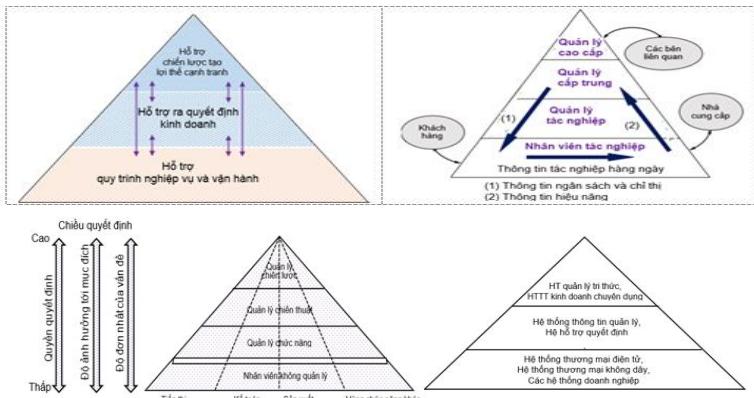


1



2

## CÁC HTTT DOANH NGHIỆP



### HTTT trong các tổ chức kinh doanh theo ba mức:

**Mức trên:** Hệ thống quản lý tri thức và hệ thống thông tin kinh doanh chuyên ngành. QL chiến lược

**Mức giữa:** HT thông tin quản lý và Hệ hỗ trợ quyết định. QL chiến thuật

**Mức dưới:** Thương mại điện tử, thương mại không dây (M-commerce: Mobile-commerce) và các hệ thống doanh nghiệp. QL chức năng (tác nghiệp)

3

## Nội dung

1. Ra quyết định và giải quyết vấn đề
2. Khái quát về HTTT quản lý
3. Các HTTT quản lý chức năng
4. Khái quát về hệ hỗ trợ quyết định
5. Các thành phần của hệ hỗ trợ quyết định
6. Hệ thống hỗ trợ nhóm
7. Hệ thống hỗ trợ điều hành
8. Dẫn luận: Công ty Generals Mills, Mỹ
9. Cty dược phẩm AstraZeneca giảm thời gian ra thị trường
10. Bốn nguyên lý và mục tiêu học tập

## 1. Ra quyết định và giải quyết vấn đề

- Giới thiệu

- Mọi tổ chức cần ra quyết định hiệu quả
- Các khóa học ra quyết định
  - Nhân viên và đơn vị kinh doanh
  - Hoàn thành mục tiêu và mục đích
- HTTT hỗ trợ **giải quyết vấn đề**:
  - Giúp ra quyết định tốt hơn và tiết kiệm hơn
- Ví dụ HTTT Trung tâm Y tế ĐH Hackensack
  - Phân tích tương tác thuốc tiềm năng “thuốc – phản ứng”.
  - Chủ đề nghiên cứu “phân tích quan hệ ngữ nghĩa (“thuốc – phản ứng”) trong văn bản y sinh.
  - Cụ thể: Thuốc trầm cảm cho bệnh nhân AIDS
  - Đầu tư hàng triệu đô la cho HTTT



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

5

5

## Các kiểu vấn đề

- **Vấn đề cấu trúc được và không cấu trúc**

- **Vấn đề cấu trúc được (structured problem)**
  - Quen thuộc, đơn giản, và các yêu cầu thông tin rõ ràng.
  - “Doanh số tuần này có cao hơn tuần trước”?
  - Chia nhỏ được thành chuỗi các bước đã được xác định tốt
  - **Tương ứng với “thuật toán hóa”: Lời giải lập trình được**
- **Vấn đề không cấu trúc được (unstructured problem)**
  - Mơ hồ do thiếu thông tin
  - “Đặc trưng khách hàng mua nhiều hàng tuần này” ?
  - Không thể chia nhỏ được thành chuỗi bước được xác định tốt
  - Cần sử dụng trực giác, lý luận, và ghi nhớ
  - **Lời giải không lập trình được**



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

6

6

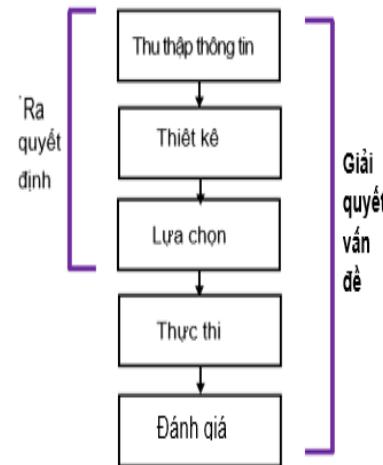
## Ra quyết định: thành phần của giải vấn đề

### • Giải vấn đề

- Hoạt động quan trọng của mọi tổ chức kinh doanh
- Người “giải vấn đề thực sự”

### • Mô hình giải vấn đề

- Mô hình ra quyết định Herbert Simon
  - Nổi tiếng
  - Ba giai đoạn: thu thập thông tin (intelligence), thiết kế (design), chọn lựa (choice)
- Mô hình giải vấn đề
  - George Huber mở rộng mô hình trên
  - Thi hành (implementation), Giám sát (monitoring) kết quả giải vấn đề



## Các giai đoạn giải vấn đề: ra quyết định

### • Ví dụ: muốn bán vải thiều Mai Siu tại Hà Nội

### • Thu thập thông tin

- Nhận dạng và xác định vấn đề hoặc cơ hội tiềm năng
- Điều tra tài nguyên và ràng buộc môi trường
  - Vấn đề: vải thiều dễ hỏng;
  - Cơ hội: giá bán buôn vải ở Hà Nội cao

### • Thiết kế

- Các giải pháp thay thế nhau (nên vài ba giải pháp) **“Mọi mô hình đều sai và có một vài mô hình dùng được”**
  - Thuê ô tô riêng / đi ô tô khách / đi bằng xe máy
  - Thời gian: lộ trình ? Chi phí ?

### • Chọn lựa

- Chọn giải pháp khả thi nhất từ các giải pháp thay thế nhau
  - Thuê ô tô riêng/đi ô tô khách/đi bằng xe máy

## Hai giai đoạn thi hành quyết định

- Ví dụ: muôn bán vải thiều Mai Siu tại Hà Nội
- Thực thi
  - Thực thi giải pháp đã lựa chọn (vận chuyển vải bằng xe máy)
  - Thông báo khách hàng, vận chuyển vải, giao trả vải, nhận tiền
- Giám sát
  - Có thông tin kết quả thực thi: thông tin phản hồi
  - Người ra quyết định tốt đánh giá giải pháp được chọn
  - Thông tin phản hồi → Điều chỉnh giải pháp được chọn
    - ví dụ, điều chỉnh lịch trình vận chuyển, cách đặt vải thiều vào sọt...
  - Thay đổi giải pháp: chọn giải pháp thay thế phù hợp



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

9

9

## Quyết định lập trình được

- Chọn lựa: nhiều nhân tố tác động đến chọn giải pháp
- Một nhân tố: quyết định lập trình được hay không
- Quyết định lập trình được
  - Có được với một quy tắc/thủ tục/phương pháp định lượng
    - Ví dụ: “hàng trong kho dưới 100 đơn vị thì cần được đặt hàng” là quyết định lập trình được vì tuân theo một quy tắc
  - Dễ dàng tin học hóa khi dùng HTTT truyền thống
    - Dễ lập trình khi số hàng trong kho  $\leq 100$  đơn vị thì đặt hàng
  - Mỗi quan hệ giữa các thành phần trong HT là cố định
  - Một dạng QĐ lập trình được: cung cấp báo cáo vấn đề thường xuyên mà mối quan hệ được xác định
- Giải pháp
  - Hầu hết quá trình tự động hóa ở hệ thống hoạch định tài nguyên doanh nghiệp / Hệ thống xử lý giao dịch
  - HTTT quản lý: các báo cáo mức cao hơn



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

10

10

## Quyết định không lập trình được

### • Tình huống

- Các tình huống bát thường hoặc đặc thù
  - Xác định chương trình đào tạo cho một nhân viên mới
  - Quyết định phát triển một sản phẩm/dịch vụ mới
  - cản nhắc lợi ích và hạn chế lắp đặt một hệ thống kiểm soát ô nhiễm nâng cấp
  - Hệ thống thông tin soạn thảo, thi hành luật
- Quyết định rất khó định lượng
- Quyết định có tính độc đáo
- Không áp dụng được các quy tắc, thủ tục chuẩn

### • Giải pháp

- Hệ hỗ trợ quyết định



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

11

11

## Tiếp cận tối ưu hóa

### • Mô hình tối ưu hóa

- Hệ HT quyết định tin học hóa là tối ưu hoặc đáp ứng
- Một quá trình tìm giải pháp tốt nhất (một trong các tốt nhất) giúp tổ chức đạt được mục tiêu của mình.
- mô hình tối ưu hóa tìm thấy giải pháp tốt nhất
  - điều kiện và giả định nhất định cho trước
  - sử dụng ràng buộc vấn đề
  - một mô hình tối ưu hóa tìm thấy lượng sản phẩm thích hợp mà tổ chức phải sản xuất để đáp ứng mục tiêu lợi nhuận



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

12

12

## Tiếp cận đáp ứng

### • Khái niệm

- Mô hình đáp ứng tìm được giải pháp tốt song không phải là tốt nhất
  - Có phương pháp tìm được giải pháp tốt
  - Rất khó đánh giá tốt nhất ?
- Không xem xét được mọi khả năng mà xem một vài khả năng tốt
- Ví dụ
  - Lựa chọn vị trí đặt cửa hàng
  - Tốt nhất: xem mọi tình huống nhưng không thể
  - Đáp ứng: khoang vùng được tốt/tốt nhất rồi mới tìm kiếm
- “Đáp ứng” là phương pháp mô hình hóa thay thế tốt
  - Quá đắt để phân tích mọi lựa chọn để lựa chọn tốt nhất,



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

13

13

## Heristic (tự khám phá) và phản hồi

### • Heuristics

- Quy tắc ngón tay cái: rules of thumb
- Chấp nhận hướng dẫn/thủ tục tìm giải pháp tốt
- Thường dùng khi ra quyết định
  - Đặt hàng trước 4 tháng khi số hàng  $\leq 20$
- Phân mềm chống thư rác
  - Heuristic dựa theo luật/phân lớp
  - Tìm ra thư có khả năng nhất
  - Không tìm tất cả các thư

### • Cảm nhận và phản hồi

- Sense and Respond
- xác định vấn đề/cơ hội (cảm nhận) và phát triển hệ thống để giải vấn đề /tận dụng cơ hội (phản hồi)
- SAR thường đòi hỏi phải tổ chức linh hoạt để thay thế dòng truyền thông

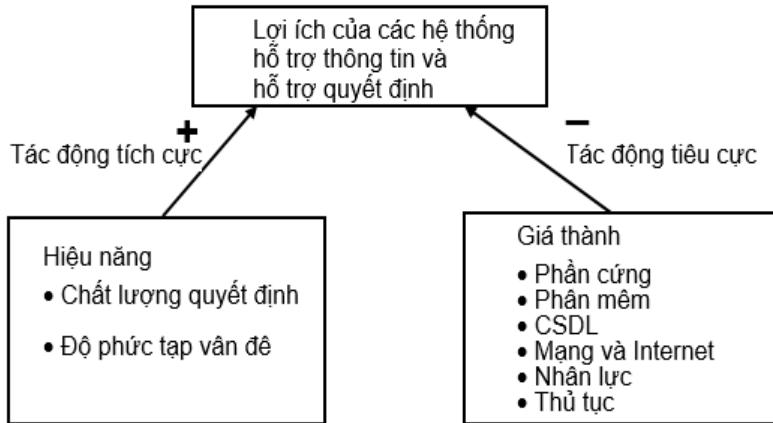


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

14

14

## Lợi ích HTTT QL và HTHTQĐ



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

15

15

## 2. Tổng quan về HTTT quản lý

### • Khái niệm

- Khái niệm
  - con người, thủ tục, cơ sở dữ liệu, và các thiết bị
  - cung cấp thông tin cho các nhà quản lý ra quyết định
  - giúp đạt được mục tiêu của tổ chức
- Lợi thế cạnh tranh: thông tin chính xác, đúng người, đúng lúc

### • Các khía cạnh

- Ngắn hạn: các báo cáo phản hồi hoạt động hàng ngày
- Mọi cấp trong toàn tổ chức [The Truth About Your Future](#)

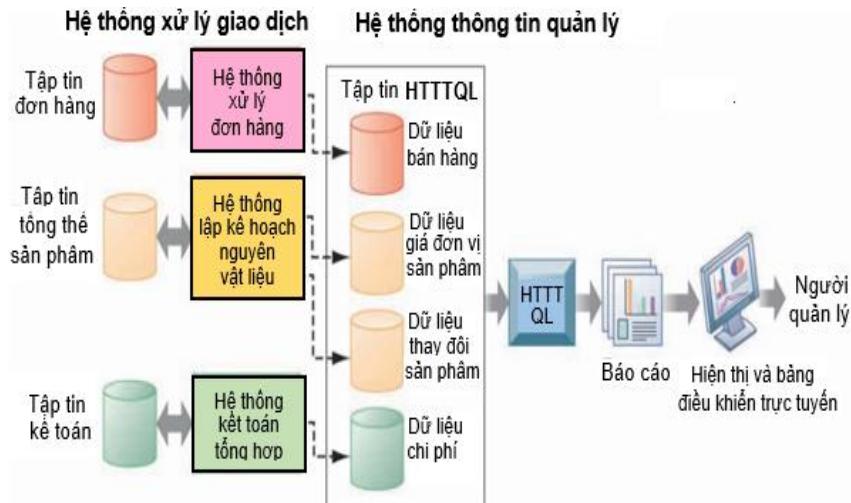


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

16

16

## HTTT quản lý

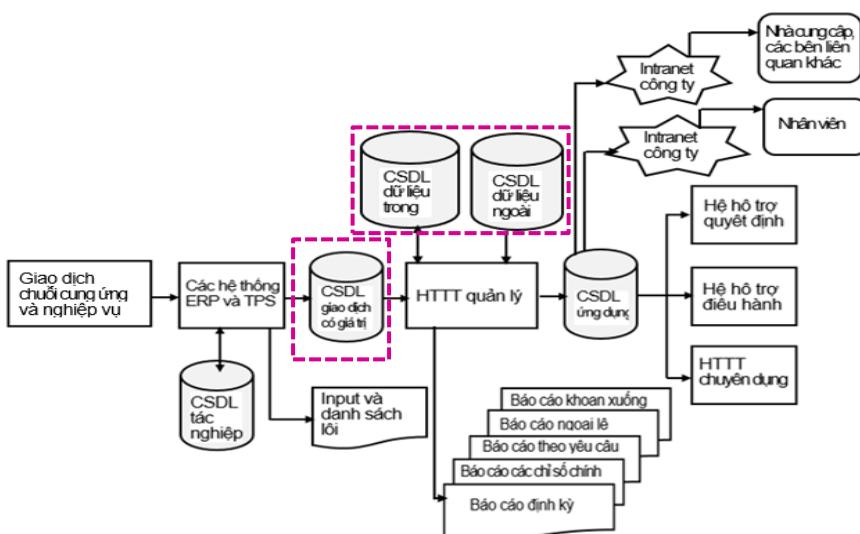


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

17

17

## HTTT quản lý: Đầu vào



HTTTQL chỉ là một trong nhiều nguồn TT quản lý (Hệ HTQD, hệ HTĐH và hệ CG cũng hỗ trợ việc ra quyết định). Tự giao dịch chuỗi cung ứng và kinh doanh

18

18

## HTTT QL: Đầu vào

- bên trong và bên ngoài, bao gồm chuỗi cung ứng (supply chain)
- được xử lý thành báo cáo dễ sử dụng cho các nhà quản lý

### • Bên trong

- Nguồn quan trọng nhất các TPS, HT ERP và CSDL liên quan
- Kho DL, kho DL chuyên (Data mart): thông tin kinh doanh giá trị, thông minh kinh doanh
- DL từ các khu vực chức năng khác

### • Bên ngoài

- DL về khách hàng, nhà cung cấp, đối thủ cạnh tranh, cỗ đông, DL khác (Internet)
- Nhóm doanh nghiệp kết nối với nhau trao đổi DL

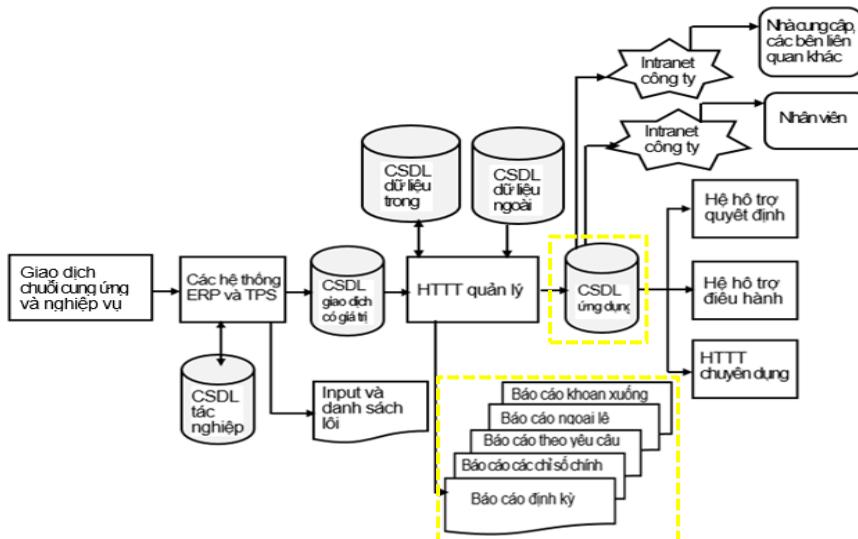


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

19

19

## HTTT quản lý: Đầu ra



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

20

20

## HTTT QL: Đầu ra

- tập kiểu các báo cáo: cung cấp từng người kịp thời
- định kỳ (scheduled), chỉ số chính (key indicator), theo yêu cầu (demand), ngoại lệ (exception), khoan xuống (drill down)

### • Các loại báo cáo

- Báo cáo định kỳ: được tạo ra theo định kỳ, hoặc theo một lịch trình, chẳng hạn như hàng ngày, hàng tuần, hoặc hàng tháng.
- Báo cáo chỉ số chính: tóm tắt các hoạt động quan trọng của ngày trước đó, sẵn sàng đầu ngày làm việc của nhà quản lý- điều hành.
- Báo cáo theo yêu cầu: báo cáo được tạo ra để cung cấp thông tin nào đó theo yêu cầu của một người (điều hành, nhà cung cấp, khách hàng).
- Báo cáo ngoại lệ: được tự động tạo ra khi một tình huống bất thường hoặc theo đòi hỏi hành động quản lý (người quản lý đặt tham số để tạo một báo cáo về mọi sản phẩm tồn kho ít hơn lượng năm ngày bán hàng hiện hành)
- Báo cáo khoan xuống (chi tiết hóa) cung cấp dữ liệu chi tiết hơn về một tình huống. Khoan xuống là một kỹ thuật trong kho dữ liệu

## Đặc điểm của HTTSQL

- Báo cáo định kỳ , chỉ số chính, theo yêu cầu, ngoại lệ, và khoan xuống giúp các nhà quản lý và điều hành ra quyết định tốt hơn, kịp thời hơn.

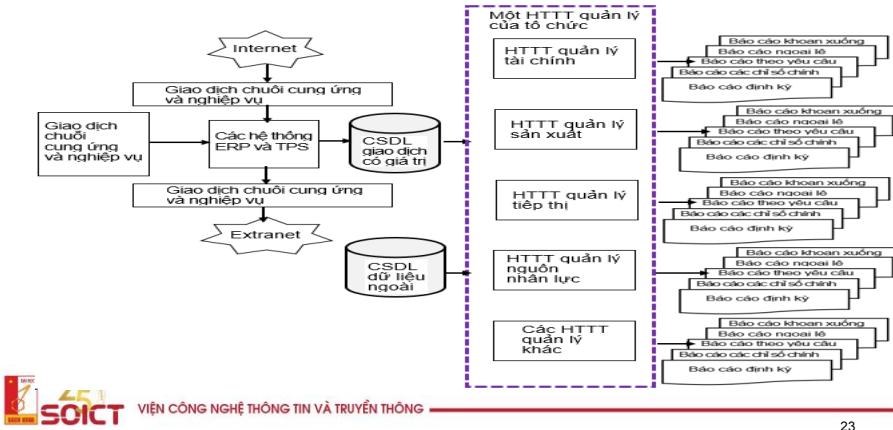
### • Các đặc điểm HTTSQL

- Cung cấp các báo cáo với các định dạng cố định và chuẩn
- Tạo ra các báo cáo bản cứng và bản mềm
- Dùng dữ liệu nội bộ được lưu trong hệ thống máy tính
- Cho phép người dùng xây dựng báo cáo của riêng họ
- Phụ thuộc yêu cầu người dùng tới các báo cáo được nhân viên hệ thống phát triển

### 3. Khía cạnh chức năng của các HTTTQL

#### • Đặt vấn đề

- Tổ chức được cấu trúc theo tuyến hoặc vùng chức năng: phân cấp theo vai trò hoặc vị trí
- MIS theo vùng chức năng truyền thống: tài chính, sản xuất, tiếp thị, nguồn nhân lực, khác



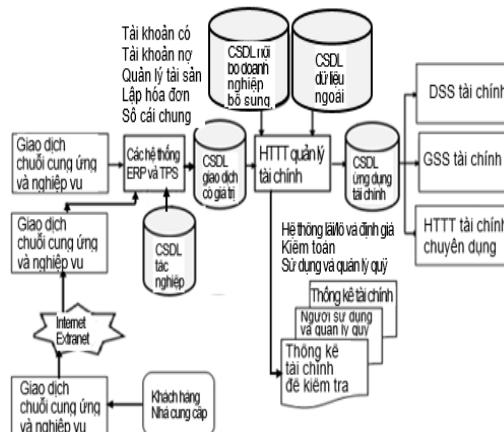
23

23

### HTTT quản lý tài chính

#### • Sơ bộ

- Phù hợp HTTTQL chung
- Yêu cầu khách hàng
- Các báo cáo
- Hệ thống lợi nhuận/ chi phí, giá
- Kiểm toán
- Dùng và quản lý quỹ
- DSS TC, GSS TC, HTTTTC chuyên sâu



#### • HTTT quản lý tài chính

- Cung cấp TT tài chính không chỉ người điều hành mà cho tập rộng rãi người cần ra quyết định tốt hơn hàng ngày (đầu tư cổ phiếu...)

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

24

24

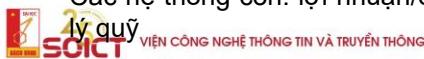
## HTTT quản lý tài chính: chức năng

### ● Chức năng

- Tích hợp thông tin tài chính và hoạt động từ nhiều nguồn (cả Internet) vào một hệ thống duy nhất
- Cung cấp sự dễ dàng truy cập dữ liệu cho người sử dụng cả về tài chính và phi tài chính, thường dùng mạng nội bộ công ty để truy cập các trang web của công ty dữ liệu và thông tin tài chính
- Tạo sự sẵn có tức thời dữ liệu tài chính để rút ngắn thời gian chu kỳ phân tích
- Cho phép phân tích dữ liệu tài chính theo nhiều chiều: thời gian, địa lý, sản phẩm, nhà máy, và khách hàng
- Phân tích hoạt động tài chính lịch sử và hiện tại
- Theo dõi và kiểm soát việc sử dụng quỹ theo thời gian

### ● Cấu trúc điển hình

- đầu vào, các hệ thống con đầu ra: xem hình vẽ trang trước
- Các hệ thống con: lợi nhuận/chi phí và giá, kiểm toán, sử dụng-quản lý quỹ



25

25

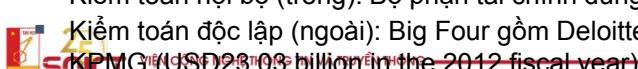
## HTTT quản lý tài chính: các hệ thống con

### ● HT con lãi/chi phí và định giá

- Các trung tâm lợi nhuận ở nhiều bộ phận: hướng tới sinh lợi nhuận. Ví dụ: phòng đầu tư của một Cty bảo hiểm và thẻ tín dụng lớn
- Các trung tâm doanh số: hướng tới tăng doanh số như các bộ phận bán hàng, tiếp thị
- Các trung tâm giá: hướng tới giảm giá như bộ phận sản xuất, R&D
- Htcon lợi nhuận, chi phí, doanh số, định giá

### ● HT con kiểm toán

- quá trình có tính hệ thống, độc lập và được làm tài liệu để có bằng chứng kiểm toán (hồ sơ, báo cáo về sự kiện/thông tin khác có liên quan và kiểm chứng được) và đánh giá nó một cách khách quan nhằm xác định mức độ các tiêu chuẩn kiểm toán (tập các chính sách, thủ tục hoặc yêu cầu) được đáp ứng. Phân tích điều kiện tài chính
- Xác định thông báo và báo cáo tài chính do MIS tạo ra có chính xác
- Kiểm toán nội bộ (trong): Bộ phận tài chính dùng nhân viên



26

26

## HTTT quản lý tài chính: HT con quỹ

### ● Dùng quỹ nội bộ

- mua hàng tồn kho bổ sung,
- nâng cấp nhà máy và thiết bị,
- thuê nhân viên mới,
- mua lại các công ty khác, mua hệ thống máy tính mới,
- tăng tiếp thị và quảng cáo, mua nguyên liệu hoặc đất đai,
- đầu tư vào sản phẩm mới, và tăng nghiên cứu và phát triển

### ● Dùng quỹ ra ngoài

- Đầu tư kinh phí dư thừa ra bên ngoài
- Tài khoản ngân hàng, cổ phiếu, trái phiếu, tín phiếu, danh tiếng, cho tương lai, quyền mua bán, ngoại tệ
- Tạo khoản vay trên Internet:



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

27

27

## HTTT quản lý sản xuất

### • Giới thiệu

- HTTTQL sản xuất: “cách mạng hóa” sản xuất bằng HTTTQL
- Cải thiện đáng kể nhiều hoạt động sản xuất
- Nhấn mạnh chất lượng và năng suất hơn → một quá trình sản xuất hiệu quả ngày càng trở nên quan trọng hơn
- Tin học hóa: từ quầy hàng tới lãnh đạo cấp cao.
- Càng nhiều công ty gia công phần mềm quá trình sản xuất

### • Nội dung sơ bộ

- HTTTQL sản xuất và kết quả ra để theo dõi & quản lý dòng chảy vật liệu, sản phẩm xuyên qua công ty.
- HTTTQL SX mọi khâu chuyển đổi nguyên vật liệu tới thành phẩm
- Công nghệ mới (chip): tạo dễ dàng dòng chuyển này
- Thành công của tổ chức phụ thuộc vào chức năng sản xuất
- Nhiều HTTTQL con.

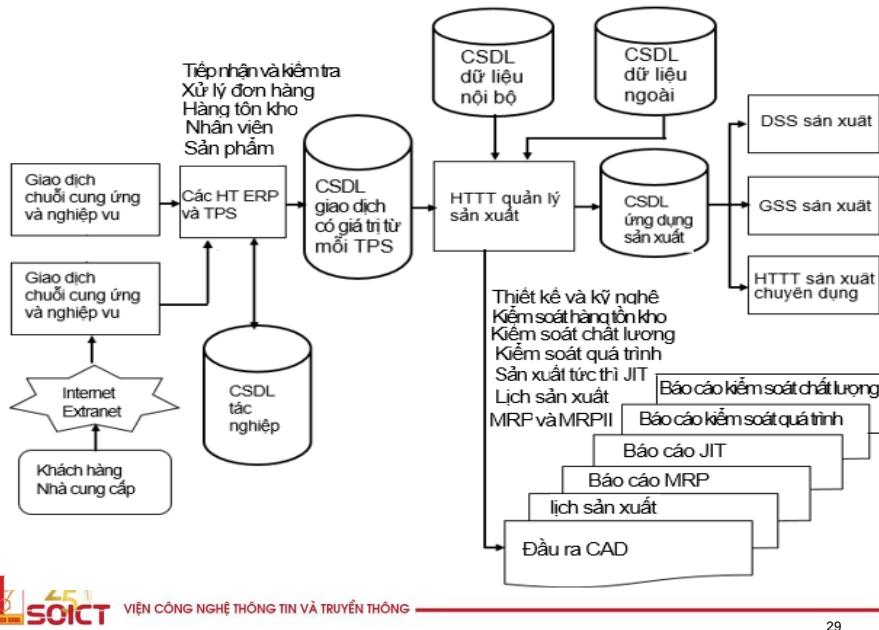


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

28

28

## HTTTQL sản xuất: sơ đồ chung



29

## HTTTQL sản xuất: Các HT con

### ● Thiết kế và kỹ nghệ

- HTTT thiết kế nhờ MT (computer-aided design: CAD) sản phẩm mới hoặc hiện có
- Dùng CAD phát triển và thiết kế sản phẩm/cấu trúc phức tạp
- Ví dụ Boeing

### ● Điều khiển lịch sản xuất và kiểm soát hàng tồn kho

- Lập kế hoạch SX và kiểm soát hàng tồn kho: rất quan trọng cho mọi công ty sản xuất
- Mục tiêu của điều khiển lập lịch: cung cấp KH chi tiết cả lịch biểu ngắn hạn và dài hạn của các CSSX
- Nhiều kỹ thuật: giảm thiểu chi phí hàng tồn kho. Khi nào và bao nhiêu hàng tồn kho cần đặt hàng? Câu trả lời reorder point: ROP.

30

## HTTSQL sản xuất: HT con điều khiển QT

- **Khái niệm**

- Theo dõi và sắp xếp dòng quá trình SX (KPQT)
- Trực tiếp điều khiển thiết bị SX: SX có máy tính hỗ trợ (computer-assisted manufacturing: CAM). Hệ thống CAM điều khiển máy khoan, dây chuyền lắp ráp, và nhiều ứng dụng khác
- Sản xuất tích hợp máy tính (Computer-integrated manufacturing: CIM): dùng máy tính liên kết các thành phần SX.  
Mục tiêu của CIM: kết hợp chặt chẽ mọi khía cạnh sản xuất, bao gồm xử lý đơn hàng, thiết kế sản phẩm, sản xuất, kiểm tra, kiểm soát chất lượng, và vận chuyển.
- Hệ thống SX linh hoạt (flexible manufacturing system: FMS): cho phép cơ sở sản xuất dễ dàng thay đổi từ một sản phẩm sang sản phẩm khác. Thay đổi quá trình

## HTTSQL sản xuất: Các HT con

- **Điều khiển quy trình**

- Phương pháp lượng đặt hàng kinh tế (economic order quantity: EOQ). Khi hàng nhỏ hơn ngưỡng (reorder point: ROP)
- Kỹ thuật lập kế hoạch yêu cầu vật liệu (MRP): tập kỹ thuật kiểm soát hàng tồn kho giúp phối hợp hàng ngàn mặt hàng tồn kho khi nhu cầu của một mặt hàng phụ thuộc vào nhu cầu khác
- Kỹ thuật hàng tồn kho và sản xuất Just-in-time (JIT): Duy trì hàng tồn kho ở mức thấp nhất mà không mất tính sẵn có của SP hoàn chỉnh.
- JIT: hết hàng tồn kho khi có đợt mua hàng nhiều

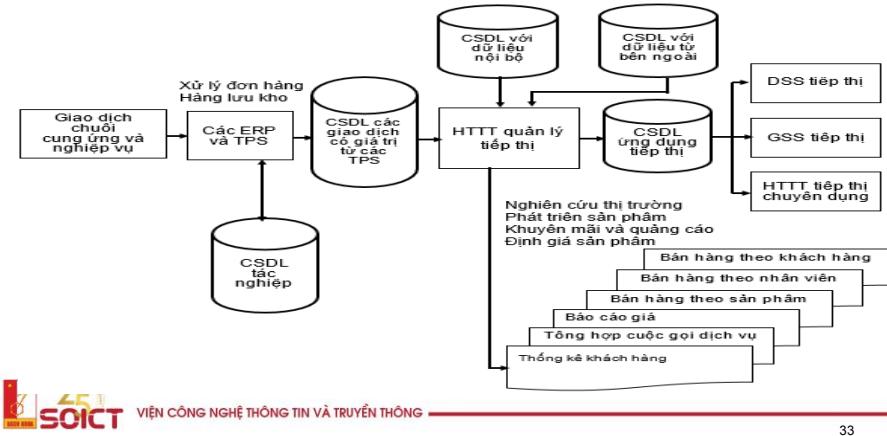
- **Kiểm soát chất lượng và kiểm thử**

- Kiểm soát chất lượng: sản phẩm đáp ứng yêu cầu người dùng ?

## HTTT quản lý tiếp thị

### • Khái niệm

- hỗ trợ quản lý phát triển sản phẩm, phân phối, quyết định giá cả, hiệu quả quảng cáo, và dự báo bán hàng.
- Xu thế sử dụng trên Internet
- Các HTTT con: nghiên cứu thị trường, phát triển sản phẩm, khuyến mãi và quảng cáo, và giá cả sản phẩm



33

## HTTT quản lý tiếp thị

### • HT con nghiên cứu thị trường

- Tiến hành một nghiên cứu chính quy về thị trường và sở thích của khách hàng
- Chương trình máy tính giúp tiến hành và phân tích kết quả điều tra, bảng hỏi, nghiên cứu thí điểm, và các cuộc phỏng vấn

### • Phát triển sản phẩm

- chuyển đổi nguyên liệu vào hàng hóa và dịch vụ hoàn thiện: tập trung chủ yếu vào các thuộc tính vật chất của sản phẩm
- Các yếu tố: năng suất máy móc, kỹ năng lao động, các yếu tố kỹ thuật, và các tài liệu
- chương trình máy tính phân tích các yếu tố khác nhau và lựa chọn sự pha trộn thích hợp của lao động, vật tư, máy móc thiết bị, và thiết kế kỹ thuật

34

## N/C thị trường: sản phẩm/thị trường

Sản phẩm mới	Nghiên cứu thị trường có thể chỉ ra khả năng tiếp nhận sản phẩm mới	Nghiên cứu thị trường có thể chỉ ra nhu cầu chưa được đáp ứng và cung cấp hiểu biết về các thị trường mới
Sản phẩm hiện có	Nghiên cứu thị trường có thể đo lường sự hài lòng khách hàng để tìm được cách thức duy trì lợi thế cạnh tranh	Nghiên cứu thị trường có thể tìm ra các vững lanh thổ mới cho các sản phẩm hoặc dịch vụ

Thị trường hiện thời

Thị trường mới

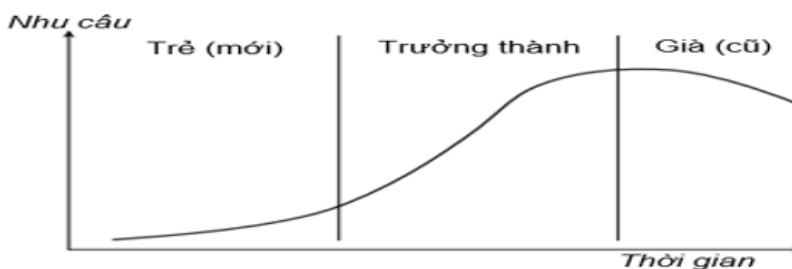


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

35

35

## N/C thị trường: nhu cầu/thời gian



Nghiên cứu thị trường khám phá các nhu cầu chưa được đáp ứng cho sản phẩm mới và giúp ước tính nhu cầu có thể. Nó có thể được sử dụng để định giá và định hình các đặc điểm kỹ thuật của sản phẩm

Nghiên cứu thị trường chỉ ra cách xây dựng thương hiệu và cạnh tranh. Nghiên cứu sự hài lòng của khách hàng chỉ ra các điểm mạnh cần được xây dựng, các điểm yếu cần được sửa chữa

Nghiên cứu thị trường chỉ ra khả năng trẻ hóa sản phẩm, có thể bằng cách kết hợp các tính năng mới hoặc tìm kiếm các thị trường mới



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

36

36

## HTTT quản lý tiếp thị

### • Xúc tiến và quảng cáo

- chức năng quan trọng nhất của mọi nỗ lực tiếp thị
- sử dụng Internet để quảng cáo và bán sản phẩm và dịch vụ
- quảng cáo truyền hình và Internet

### • Định giá SP

- một chức năng tiếp thị quan trọng và phức tạp: giá bán lẻ, giá bán buôn, giảm giá
- phát triển chính sách giá cả để tối đa tổng doanh thu bán hàng: chuong trình máy tính phân tích mối quan hệ giữa giá và tổng doanh thu

### • Phân tích bán hàng

- rất quan trọng để xác định sự đóng góp của sản phẩm, nhân viên bán hàng, và khách hàng đóng góp vào lợi nhuận
- Một số báo cáo được tạo ra giúp đưa ra quyết định bán hàng tốt
  - Báo cáo bán hàng theo sản phẩm chính
  - Báo cáo bán hàng theo nhân viên bán hàng
  - Báo cáo bán hàng theo khách hàng



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

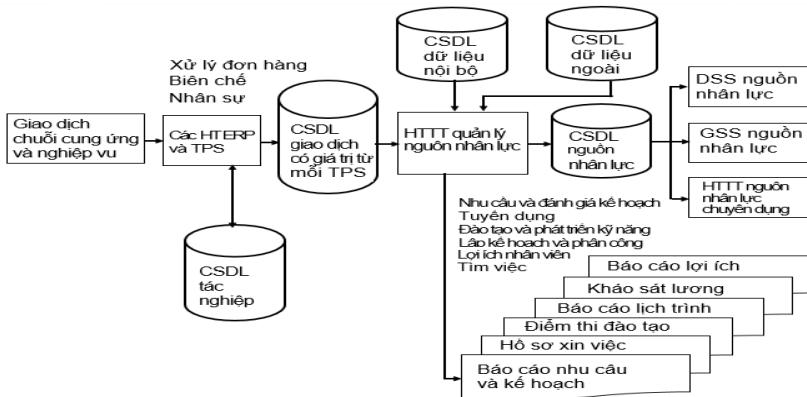
37

37

## HTTT quản lý nhân lực

### • Khái niệm

- quan tâm tới các hoạt động liên quan đến người lao động và người lao động tiềm năng
- Thời gian: quá khứ, hiện tại và tương lai
- Đóng vai trò quan trọng đảm bảo thành công của tổ chức



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

38

38

## HTTT quản lý nhân lực

### • Một số nội dung

- Lập kế hoạch nguồn nhân lực: khía cạnh đầu tiên xác định nhân viên và nhu cầu nhân lực
- Lựa chọn và tuyển dụng: chương trình máy tính sắp xếp các nỗ lực tuyển dụng, kiểm tra kỹ năng nhân viên tiềm năng, sàng lọc người xin việc qua Internet
- Đào tạo và phát triển kỹ năng: yêu cầu đào tạo rất cụ thể cho nhân viên mới như văn hóa tổ chức, định hướng, tiêu chuẩn, và mong đợi của tổ chức.
- Lập kế hoạch và sắp xếp công việc
- Tiền lương và quản lý tiền lương
- Tìm việc. Nhân viên rời khỏi công ty vì nhiều lý do. Công ty cung cấp dịch vụ tìm việc giúp nhân viên thực hiện quá trình chuyên đổi.

## HTTT quản lý khác

### • HTTT QL kế toán

- Liên quan tới MIS tài chính
- thực hiện một số hoạt động quan trọng, cung cấp thông tin tổng hợp về các khoản phải trả, các khoản phải thu, biên chế, và nhiều ứng dụng khác

### • Hệ thống thông tin địa lý (GIS)

- Geographic Information Systems
- Trực quan hóa dữ liệu dưới dạng đồ họa
- hệ thống máy tính có khả năng lắp ráp, lưu trữ, thao tác, và hiển thị thông tin địa lý tham chiếu: dữ liệu xác định theo vị trí của nó

## 4. Khái quát về HHTQĐ

- Hệ HTQĐ nhiều đặc trưng  $\Rightarrow$  trở thành công cụ hỗ trợ QL hiệu quả

### • Đặc trưng của HHTQĐ

- Cung cấp truy cập nhanh đến thông tin.
- Xử lý một lượng lớn dữ liệu từ các nguồn khác nhau
- Cung cấp và trình bày linh hoạt các báo cáo
- Cung cấp cả hai định hướng văn bản và đồ họa
- Hỗ trợ phân tích khoan sâu
- Thực hiện phân tích và so sánh phức tạp, tinh vi sử dụng các gói phần mềm tiên tiến
- Hỗ trợ tối ưu hóa, pháp thoả mãn, và cách tiếp cận heuristic
- Thực hiện phân tích mô phỏng: khả năng của HHTQĐ để sao chép các tính năng của một hệ thống thực, nơi có liên quan tới tính khả năng hoặc tính không chắc chắn

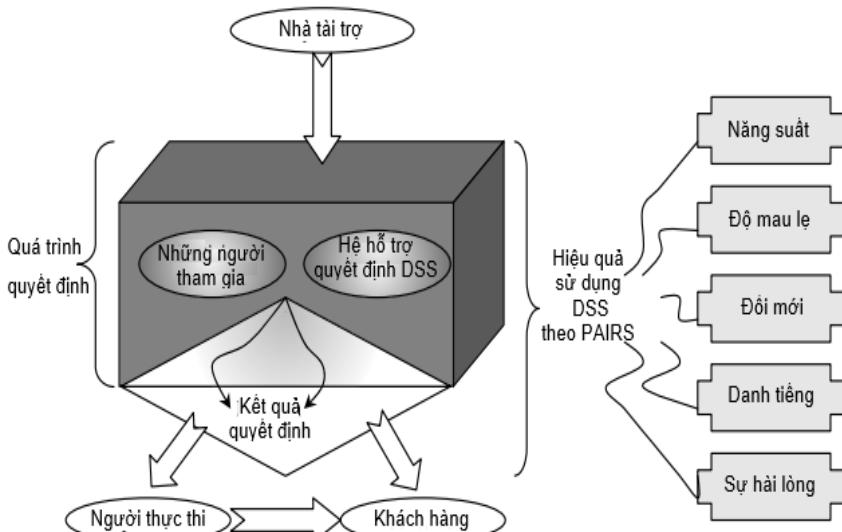


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

41

41

## Vai trò của HHTQĐ



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

42

42

## HHTQĐ: Một số ví dụ

Tổ chức hoặc ứng dụng	Mô tả
ING Direct	Công ty dịch vụ tài chính dùng DSS để tổng hợp hiệu năng tài chính của một ngân hàng. Ngân hàng cần có cơ chế đo lường và theo dõi để xác định nó thành công ra sao và cách cải tiến để lên kế hoạch thời gian thực
Cinergy Corporation	Tiên ích số hóa phát triển một HHTQĐ để rút gọn thời gian và nỗ lực đòi hỏi ra quyết định
U.S. Army	Lực lượng quân đội Mỹ phát triển một HHTQĐ để trợ giúp tân binh, huấn luyện, và đào tạo nhập ngũ quân đội. HHTQĐ này sử dụng mô phỏng kết hợp các đặc trưng What-If.
National Audubon Society	Tổ chức phát triển một HHTQĐ được gọi là Energy Plan (EPLAN) để phân tích ảnh hưởng của chính sách năng lượng Mỹ tới môi trường
Hewlett-Packard	Công ty máy tính phát triển một HHTQĐ được gọi là Quality Decision Management để hỗ trợ nâng cao chất lượng sản phẩm và dịch vụ
State of Virginia	Bang Virginia phát triển HHTQĐ đánh giá vận chuyển (the Transportation Evacuation Decision Support System: TEDSS) để xác định con đường tốt nhất để người dân tránh được tai họa hạt nhân từ các nhà máy điện hạt nhân.

## Năng lực của HHTQĐ

- Giới thiệu
  - Phát triển HHTQĐ nhằm mục đích:
    - Linh hoạt hơn HTTSQL
    - Tăng cường năng lực hỗ trợ ra quyết định trong nhiều tình huống
  - HHTQĐ hỗ trợ giải vấn đề
    - toàn bộ/hầu hết các khâu
    - Kiểu thường xuyên quyết định
    - Kiểu cấu trúc vấn đề
  - Đối tượng sử dụng HHTQĐ: người ra quyết định ở mọi mức
- Tiếp cận HHTQĐ
  - Hỗ trợ mọi mức của quá trình ra quyết định
  - HHTQĐ chỉ thi hành một vài năng lực
  - Mục đích và phạm vi sử dụng một HHTQĐ quy định tập con năng lực của HHTQĐ đó

## HHTQĐ: Tập năng lực

- **Hỗ trợ các giai đoạn giải vấn đề**

- Hỗ trợ một vài pha trong thu thập TT, thiết kế, lựa chọn, thực thi, giám sát
- Hỗ trợ nhiều tiếp cận trong mỗi pha: linh hoạt cho người ra QĐ

- **Hỗ trợ ra quyết định thường xuyên khác nhau**

- Từ QĐ đơn nhất (one-of-a-kind) tới QĐ được lặp đi lặp lại
- QĐ đơn nhất:
  - xuất hiện chỉ một vài lần trong cuộc sống của tổ chức
  - doanh nghiệp nhỏ: có thể chỉ xảy ra một lần
  - Xây dựng nhà máy tài khu vực khác trong nước
  - HHTQĐ chuyên biệt (ad hoc DSS)
- QĐ lặp đi lặp lại:
  - Vài lần hoặc hơn trong một năm
  - HHTQĐ tổ chức (institutional DSS)
  - Vài lần/năm và được tinh chỉnh theo thời gian: Vấn đề Danh mục đầu tư, quyết định đầu tư, lập lịch sản xuất
  - Vài lần/ngày: Giải vấn đề dựa vào máy tính, DSS giám sát từng giây



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

45

45

## HHTQĐ: Tập năng lực (tiếp)

- **Hỗ trợ giải vấn đề với nhiều mức cấu trúc**

- Cấu trúc được và lập trình được cao → không cấu trúc và không lập trình được
- Vấn đề cấu trúc được: các sự kiện và quan hệ đã biết → HHTQĐ đơn giản
- Vấn đề không cấu trúc và nửa cấu trúc → HHTQĐ phức tạp
  - quan hệ giữa các DL không tường minh,
  - DL ở nhiều định dạng khó thao tác
  - Yêu cầu thông tin quyết định có thể chưa biết trước
  - DSS hỗ trợ phân tích đầu tư tình vi và không cấu trúc → tạo lợi nhuận đáng kể cho thương nhân và nhà đầu tư
  - Một vài phần mềm DSS lập trình đặt lệnh mua và bán tự động



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

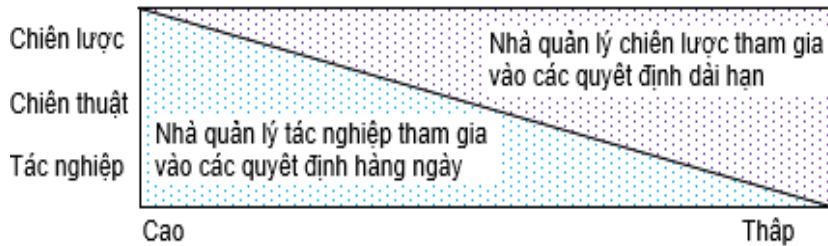
46

46

## HHTQĐ: Tập năng lực (tiếp)

### • Hỗ trợ nhiều cấp độ ra quyết định

- Trong một tổ chức: một HHTQĐ hỗ trợ nhiều cấp
- Mức quản lý tác nghiệp: ra QĐ hàng ngày và thường xuyên
- Mức QĐ chiến thuật: lập kế hoạch và kiểm soát đúng cách
- Mức QĐ chiến lược: thông tin dài hạn



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

47

47

## So sánh HHTQĐ với HTTT quản lý

<i>Yếu tố</i>	<i>HHTQĐ</i>	<i>HTTTQL</i>
<b>Loại vấn đề</b>	HHTQĐ có thể giải quyết vấn đề phi cấu trúc mà không dễ dàng lập trình	HTTTQL thường chỉ được dùng để giải quyết vấn đề cấu trúc được
<b>Kiểu người sử dụng</b>	Một HHTQĐ hỗ trợ cá nhân, nhóm nhỏ, và toàn bộ tổ chức. Người sử dụng thường có quyền kiểm soát nhiều hơn đối với HHTQĐ	HTTTQL chủ yếu hỗ trợ tổ chức. Người sử dụng ít có quyền kiểm soát HHTTQL
<b>Kiểu hỗ trợ</b>	HHTQĐ hỗ trợ mọi khía cạnh và mọi giai đoạn ra quyết định, nó không thay thế việc ra quyết định, con người phải ra quyết định	Một vài HHTTQL ra quyết định tự động và thay thế người ra quyết định.
<b>Nhân mạnh</b>	HHTQĐ nhân mạnh quyết định thực tiễn và phong cách ra quyết định	HTTTQL thường chỉ nhân mạnh tới thông tin.
<b>Tiếp cận</b>	HHTQĐ là một hệ hỗ trợ trực tiếp, cung cấp báo cáo tương tác trên màn hình máy tính	HTTTQL thường là hệ hỗ trợ gián tiếp khi dùng báo cáo ra thường xuyên.
<b>Hệ thống</b>	Thiết bị cung cấp HHTQĐ thường là trực tuyến (kết nối trực tiếp với hệ thống máy tính), liên quan thời gian thực (cung cấp kết quả ngay lập tức). Trạm cuối và màn hình hiển thị: ví dụ cung cấp thông tin trả lời tức thời cho câu hỏi	HTTTQL dùng các báo cáo được in ra, được chuyển giao cho người quản lý một lần mỗi tuần, không cung cấp kết quả ngay lập tức.
<b>Tốc độ</b>	Vì HHTQĐ là linh hoạt, được thực hiện bởi người sử dụng, nó thường mất ít thời gian hơn để phát triển và có khả năng đáp ứng yêu cầu người sử dụng tốt hơn	Thời gian đáp ứng của HHTTQL thường là dài hơn.
<b>Đầu ra</b>	Báo cáo từ HHTQĐ thường định hướng màn hình, cung năng lực tạo báo cáo tại máy in	HTTTQL thường được định hướng in ra các báo cáo và tài liệu.
<b>Phát triển</b>	Người sử dụng HHTQĐ thường liên quan trực tiếp hơn tới sự phát triển HHTQĐ. Sự tham gia của người sử dụng thường theo nghĩa hệ thống tốt hơn để cung cấp hỗ trợ tốt hơn. Với mọi HHTQĐ, sự tham gia của người sử dụng là yêu tố quan trọng nhất cho sự phát triển thành công hệ thống	HTTTQL thường có vòng đời nhiều năm và thường được phát triển cho những người không còn thực thi các công việc được hỗ trợ bởi HHTTQL đó



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

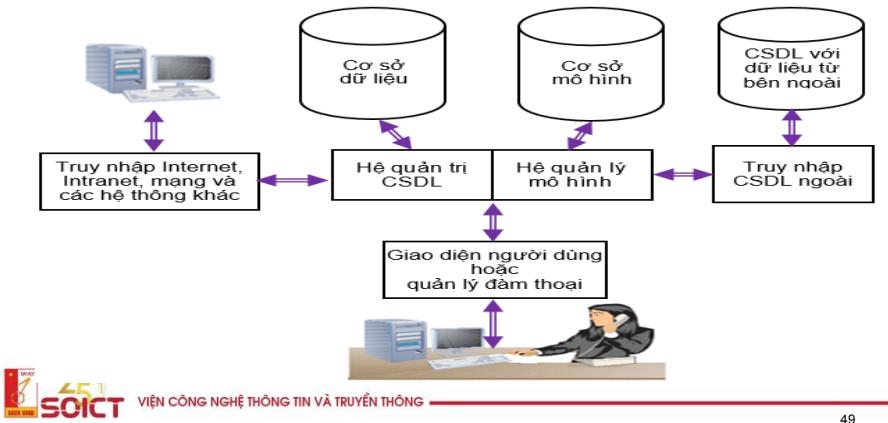
48

48

## 5. Các thành phần của HHTQĐ

- **Thành phần chung**

- Cốt lõi: CSDL và Cơ sở mô hình
- Giao diện người dùng (bộ quản lý hội thoại: dialogue manager)
- Hình vẽ: Mô hình khái niệm HHTQĐ



49

49

## HHTQĐ: Cơ sở dữ liệu

- **Đặc trưng**

- Cho phép phân tích định tính để ra quyết định
- DL nội bộ đa dạng của công ty: CSDL, kho DL, kho DL chuyên đề

- **HHTQĐ định hướng DL**

- Data-driven DSS (DDSS)
- Lấy thông tin hàng tồn kho, bán hàng, nhân viên, sản xuất, tài chính, kế toán... hỗ trợ quyết định để giảm chi phí hàng tồn kho
- Khai phá dữ liệu và thông minh kinh doanh (BI)
- DDSS y tế: bác sĩ truy cập hồ sơ y tế đầy đủ của bệnh nhân
- Lưu ý vấn đề riêng tư
- Có thể kết nối lấy DL ngoài

50

## HHTQĐ: Cơ sở mô hình

- **Đặc trưng**

- Cho phép phân tích định lượng để ra quyết định
- Dữ liệu nội bộ và bên ngoài
- Cơ sở mô hình: các mô hình miền bài toán

- **HHTQĐ định hướng mô hình**

- Model-driven DSS (MDSS)
- Quản lý mô hình cho phép người dùng truy cập nhiều mô hình, tạo kịch bản theo mô hình và trực quan hóa kết quả
- Procter & Gamble: MDSS sắp xếp hợp lý hóa dòng chảy nguyên vật liệu và sản phẩm từ các nhà cung cấp tới khách hàng: tiết kiệm được hàng trăm triệu US\$ chi phí chuỗi cung ứng.
- Tiện lợi cho dự đoán hành vi khách hàng
- LoanPerformance ([www.loanperformance.com](http://www.loanperformance.com)): dùng MDSS hỗ trợ dự báo khách hàng có thể trả hoặc vỡ nợ
- Highmark (bảo hiểm y tế): dùng MDSS để dự đoán gian lận



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

51

51

## HHTQĐ: Cơ sở mô hình (tiếp)

- **Phần mềm quản lý mô hình**

- Model management software | hệ thống quản lý mô hình
- các mô hình: tài chính, thống kê, đồ họa, quản lý dự án
- Sử dụng một hoặc phối hợp nhiều mô hình theo nhu cầu

<i>Loại mô hình</i>	<i>Mô tả</i>	<i>Phần mềm</i>
Tài chính	Cung cấp dòng tiền mặt, tỷ lệ hoàn vốn nội bộ, và phân tích đầu tư khác	Bảng tính như Microsoft Excel
Thống kê	Cung cấp thống kê tổng hợp, dự báo xu hướng, kiểm tra giả thuyết, v.v.	Chương trình thống kê như SPSS hoặc SAS
Đồ họa	Hỗ trợ ra quyết định trong thiết kế, phát triển và dựng màn hình đồ họa dữ liệu và thông tin	Chương trình đồ họa như Microsoft PowerPoint
Quản lý dự án	Xử lý & điều phối các dự án lớn, cũng dùng để xác định các hoạt động & bài toán then chốt mà có thể trì hoãn hoặc gây nguy hiểm cho toàn bộ dự án nếu không được hoàn thành kịp thời và chi phí hiệu quả	Phần mềm quản lý dự án như Microsoft Project



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

52

52

## HHTQĐ: Giao diện người dùng

- Vai trò

- user interface / dialogue manager
- Cho phép tương tác người dùng với DSS để nhận thông tin
- Hỗ trợ mọi phương diện truyền thông giữa người dùng và phần cứng & phần mềm tạo thành DSS
- Quan niệm người dùng: Giao diện chính là DSS

- Người ra quyết định mức trên

- ít quan tâm: nơi mà thông tin đến hoặc cách thông tin được thu thập
- Quan tâm nhiều: thông tin dễ hiểu và dễ truy cập
- Giao diện thân thiện, phù hợp người dùng



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

53

53

## 6. Hệ hỗ trợ làm việc nhóm

- GSS và vai trò

- Group Support System | Group Decision Support System
- DSS hỗ trợ ra quyết định ở mức cá nhân
- Đặc trưng của tổ chức: hoạt động nhóm
- Ra quyết định ở mức trên (chiến lược và chiến thuật) cần làm việc nhóm
- GSS hỗ trợ ra quyết định ở mức nhóm làm việc trên máy tính
- GSS = DSS + phần mềm hỗ trợ hiệu quả môi trường ra quyết định nhóm

- Ứng dụng

- dùng trong hầu hết ngành công nghiệp, chính phủ và quân đội
- Kiến trúc sư+kiến trúc sư và nhà xây dựng tạo kế hoạch tốt nhất và hợp đồng cạnh tranh
- Nhà sản xuất: DSS nối nhà cung cấp nguyên liệu tới HT của họ
- Mathcad Enterprise: cho phép tạo, chia sẻ, và tái sử dụng tính toán
- Có thể dùng phương tiện xã hội cho GSS

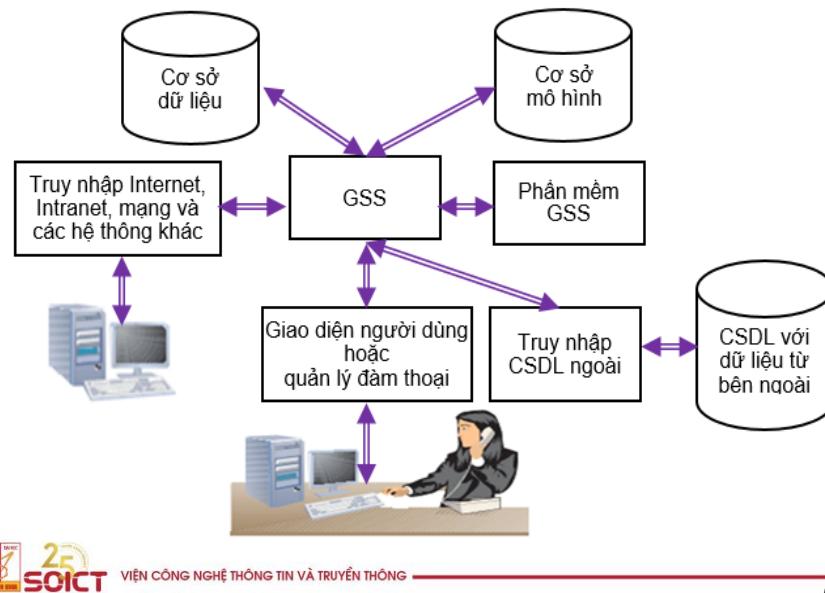


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

54

54

## Mô hình GSS = DSS + phần mềm GSS



55

55

## Đặc trưng GSS: nâng cao ra quyết định

- **Nâng cao quyết định**
  - “Hai cái đầu tốt hơn một”; “Một cây – không, ba cây – hòn núi cao”
  - Độc đáo để có quyết định tốt hơn
  - DSS hỗ trợ cá nhân + độc đáo làm việc nhóm
  - Độc đáo: có thể không là mặt đối mặt
- **Đặc trưng: Thiết kế đặc biệt**
  - Thủ tục đặc biệt với thiết bị, tiếp cận độc đáo
  - Thủ tục: thúc đẩy tư duy sáng tạo, truyền thông hiệu quả, và kỹ thuật ra quyết định nhóm tốt
- **Dễ sử dụng**
  - GSS phải dễ hiểu và dễ sử dụng
- **Linh hoạt**
  - Thành viên nhóm: phong cách và sở thích riêng. Một mặt cần nâng cao kỹ năng làm việc nhóm song GSS cũng cần linh hoạt

56

## Các đặc trưng GSS

### • Hỗ trợ ra quyết định

- hỗ trợ các phương pháp ra quyết định khác nhau: *phương pháp Delphi*, *phương pháp khởi nguồn ý tưởng* (*brainstorming*), *đồng thuận nhóm* (*group consensus approach*), *kỹ thuật nhóm danh nghĩa* (*nominal group technique*),

### • Đầu vào ẩn danh

- *anonymous*: người nạp dữ liệu được “ẩn danh” với các thành viên khác trong nhóm (không biết ai đã nạp dữ liệu).
- Đầu vào ẩn danh: người ra quyết định nhóm tập trung vào giá trị đầu vào mà không xem xét người nào đã ra mỗi quyết định
- Bỏ phiếu kín
- sử dụng đầu vào ẩn danh đưa ra quyết định tốt hơn  $\Leftrightarrow$  có thể dẫn đến tinh huống “thiêu đốt” (*flaming*), thành viên nhóm ẩn danh viết ra lời lăng mạ hoặc thậm chí những lời tục tĩu trên GSS



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

57

57

## Các đặc trưng GSS

### • Giảm thiểu hành vi nhóm tiêu cực

- thiêu đốt khi đầu vào ẩn danh: một dạng hành vi nhóm tiêu cực
- Thông đồng, cá nhân chi phối
- tư duy tập thể (*groupthink*)
- Thủ tục lập kế hoạch và quản lý các cuộc họp nhóm hiệu quả

### • Truyền thông song song và đơn nhất

- moi thành viên giải quyết vấn đề/có ý kiến cùng một lúc bằng cách nhập ý kiến từ máy tính cá nhân hoặc máy trạm
- nhận định và vấn đề được hiển thị tức thì trên máy tính cá nhân hoặc trạm làm việc
- Truyền thông song song (*parallel communication*), truyền thông hợp nhất (*unified communications*)

### • Lưu trữ hồ sơ tự động

- lưu giữ tự động hồ sơ chi tiết của cuộc họp nhóm
- bình luận, ý kiến được nhập vào từ máy tính hoặc máy trạm của một thành viên có thể được ghi ẩn danh

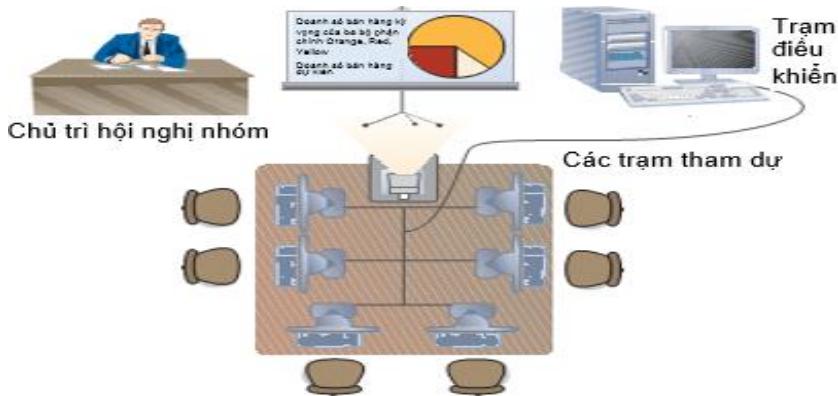


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

58

58

## Phần mềm hệ hỗ trợ nhóm



- Phần mềm GSS (*groupware/ workgroup software*).
- Lịch biểu điện tử dùng chung (*shared electronic calendars*)
- Phần mềm GSS ngày càng được tích hợp vào gói phần mềm hiện có (gói HTXLGD và hệ thống ERP)

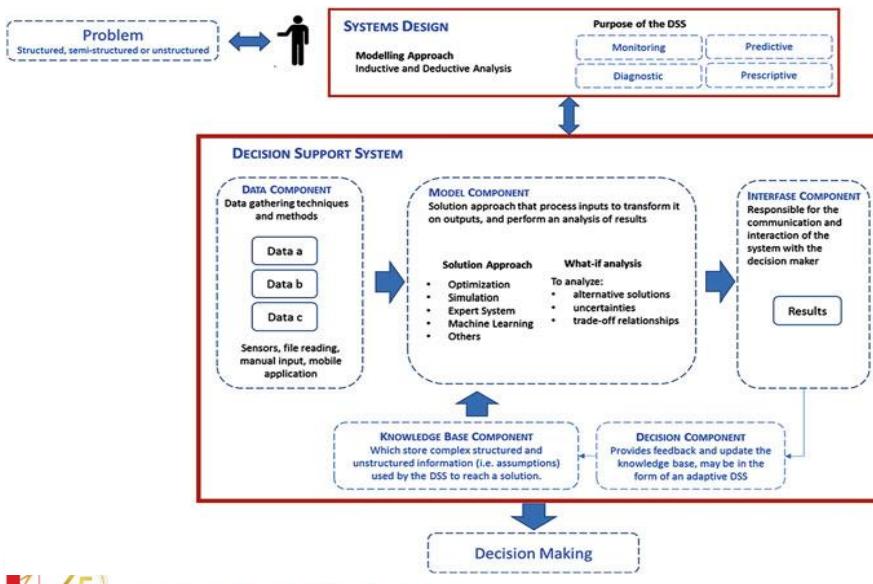


VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

59

59

## Hệ HTQĐ ví dụ



Các thành phần cốt lõi của một DSS hiện đại cần có giao diện người dùng thân thiện

60

60

## 7. HT hỗ trợ điều hành

### • Khái niệm

- *Executive Support System*: ESS
- HHTQĐ chuyên dụng hỗ trợ giám đốc điều hành cấp cao
- ESS, còn được gọi là hệ thống thông tin điều hành (*Executive Information System*: EIS), hỗ trợ hoạt động ra quyết định của các thành viên Ban giám đốc, người chịu trách nhiệm cho các cổ đông
- Các tầng như hình vẽ



## Yêu cầu và khả năng HT hỗ trợ điều hành

### • Yêu cầu

- Phù hợp với cá nhân giám đốc điều hành,
- Dễ sử dụng,
- Có khả năng khoan xa,
- Hỗ trợ nhu cầu về dữ liệu bên ngoài,
- Trợ giúp trong các tình huống có mức độ không chắc chắn cao,
- Có định hướng tương lai,
- Được liên kết với các quá trình kinh doanh giá trị gia tăng

### • Khả năng

- Hỗ trợ xác định một tầm nhìn tổng thể
- Hỗ trợ lập kế hoạch chiến lược (*Strategic planning*)
- Hỗ trợ tổ chức và nhân sự chiến lược
- Hỗ trợ kiểm soát chiến lược
- Hỗ trợ quản lý khủng hoảng



63



64