# Data models ( belongs to Programability)

## Introduction to Big Data

# Index

Universidad
Francisco de Vitoria
**UFV** Madrid

**Data models describe data characteristics:**

Structures

Operations

Constrictions

Structures

Structured

Semi structured

Unstructured

## Structured

- **Relational models:**
  - **Table colection**
  - **Without duplicates**

| ID | FName | LName | Department | Title | Salary |
|---|---|---|---|---|---|
| 202 | John | Gonzales | IT | DB Specialist | 104750 |
| 203 | Mary | Roberts | Research | Director | 175400 |
| 204 | Janaki | Rao | HR | Financial Analyst | 63850 |
| 205 | Alex | Knight | IT | Security Specialist | 123500 |
| 206 | Pamela | Ziegler | IT | Programmer | 85600 |
| 207 | Harry | Dawson | HR | Director | 115450 |
| ~~207~~ | ~~Harry~~ | ~~Dawson~~ | ~~HR~~ | ~~Director~~ | ~~115450~~ |

- **Sin "tuples" no permitidos**

| ID | Fname | Lname | Department | Title | Salary |
|---|---|---|---|---|---|
| 202 | John | Gonzales | IT | DB Specialist | 104750 |
| 203 | Mary | Roberts | Research | Director | 175400 |
| 204 | Janaki | Rao | HR | Financial Analyst | 63850 |
| 205 | Alex | Knight | IT | Security Specialist | 123500 |
| 206 | Pamela | Ziegler | IT | Programmer | 85600 |
| 207 | Harry | Dawson | HR | Director | 115450 |
| Jane | Doe | 208 | Res. Associate | 65800 | Research |

**Structured**

- **Relational model:**
  - **Primary Key**
  - **Constrictions**

## Employee

| ID: Int Primary key | Fname: string Not null | Lname: string Not null | Department: Enum (HR, IT, Research, Business) | Title: string | Salary: int > 25000 |
|---|---|---|---|---|---|
| 202 | John | Gonzales | IT | DB Specialist | 104750 |
| 203 | Mary | Roberts | Research | Director | 175400 |
| 204 | Janaki | Rao | HR | Financial Analyst | 63850 |
| 205 | Alex | Knight | IT | Security Specialist | 123500 |
| 206 | Pamela | Ziegler | IT | Programmer | 85600 |
| 207 | Harry | Dawson | HR | Director | 115450 |
| ~~Jane~~ | ~~Doe~~ | ~~208~~ | ~~Res. Associate~~ | ~~65800~~ | ~~Research~~ |

**Foreign Keys**

EmpSalaries

| EmpID | Date | Salary |
|---|---|---|
| 202 | 1/1/2016 | 104750 |
| 203 | 2/15/1016 | 175400 |
| 204 | 6/1/2015 | 63850 |
| 205 | 9/15/2015 | 123500 |
| 206 | 10/1/2015 | 85600 |
| 207 | 4/15/2015 | 115450 |
| 202 | 9/15/2014 | 101250 |
| 204 | 3/1/2015 | 48000 |
| 207 | 9/15/2013 | 106900 |
| 205 | 10/1/2014 | 113400 |

Foreign Key

EmpSalaries.EmpID References Employees.ID

Primary key

## Structured

- **Relational models:**
  - **Foreign key**
  - **Table relationships through the Foreing keys**

**Employee**

| ID: Int Primary key | Fname: string Not null | Lname: string Not null | Department: Enum (HR, IT, Research, Business) | Title: string | Salary: int > 25000 |
|---|---|---|---|---|---|
| 202 | John | Gonzales | IT | DB Specialist | 104750 |
| 203 | Mary | Roberts | Research | Director | 175400 |
| 204 | Janaki | Rao | HR | Financial Analyst | 63850 |
| 205 | Alex | Knight | IT | Security Specialist | 123500 |
| 206 | Pamela | Ziegler | IT | Programmer | 85600 |
| 207 | Harry | Dawson | HR | Director | 115450 |
| ~~Jane~~ | ~~Doe~~ | ~~208~~ | ~~Res. Associate~~ | ~~65800~~ | ~~Research~~ |

**EmpSalaries**

| EmpID | Date | Salary |
|---|---|---|
| 202 | 1/1/2016 | 104750 |
| 203 | 2/15/1016 | 175400 |
| 204 | 6/1/2015 | 63850 |
| 205 | 9/15/2015 | 123500 |
| 206 | 10/1/2015 | 85600 |
| 207 | 4/15/2015 | 115450 |
| 202 | 9/15/2014 | 101250 |
| 204 | 3/1/2015 | 48000 |
| 207 | 9/15/2013 | 106900 |
| 205 | 10/1/2014 | 113400 |

Foreign key

EmpSalaries.EmpID References Employees.ID

Primary key

Semi structured

- **Semi structured**
  - **HTML, XML, JSON,...**



## A Simple HTML Example

This is the first paragraph.

- List item 1
- List item 2

**This is a bolded text.**

```
<!DOCTYPE html>
<html>
<body>

<h1>A Simple HTML Example</h1>

<p title="undecided so far">
This is the first paragraph.
<li> List item 1 </li>
<li> List item 2 </li>
</p>

<p><b>
This is a bolded text.
</b></p>

</body>
</html>
```
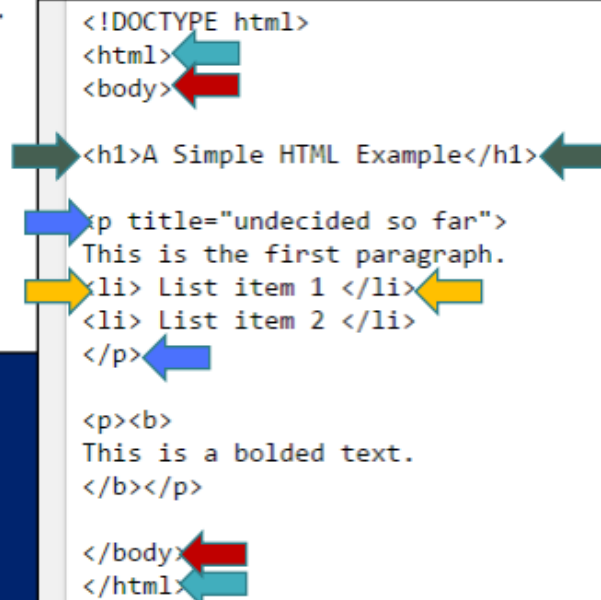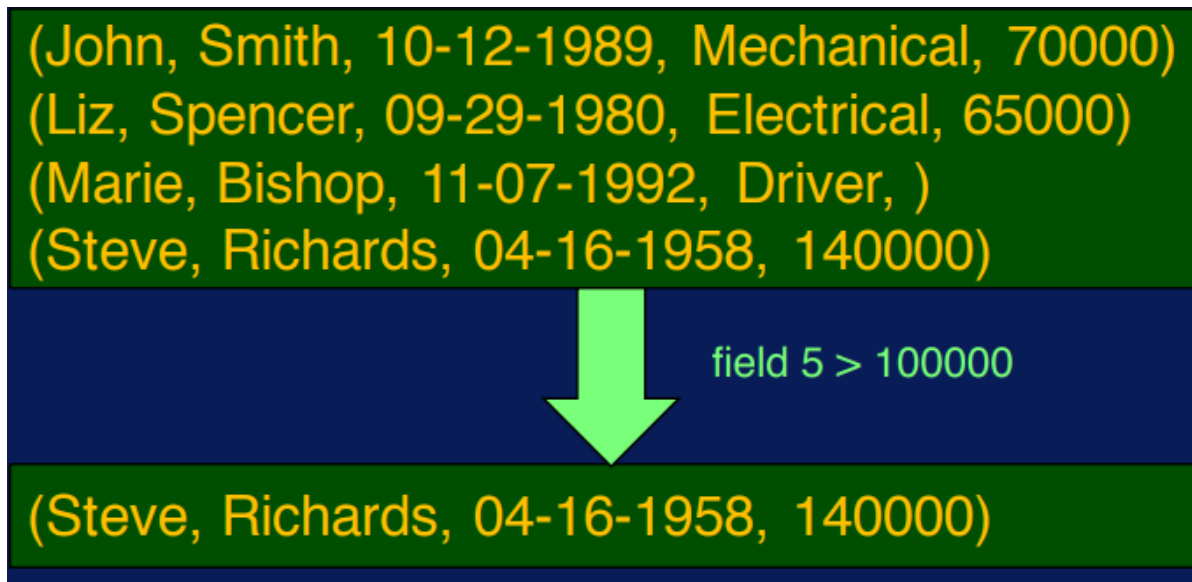
Operations

- **Sub setting:** given a data set and a condition
  - Find a subset that fulfils the condition

(John, Smith, 10-12-1989, Mechanical, 70000)
(Liz, Spencer, 09-29-1980, Electrical, 65000)
(Marie, Bishop, 11-07-1992, Driver, )
(Steve, Richards, 04-16-1958, 140000)

field 5 > 100000

(Steve, Richards, 04-16-1958, 140000)

**Operations**

- **Substructure extraction:** given a set of data with an
  - Extract a part of that structure with its elements

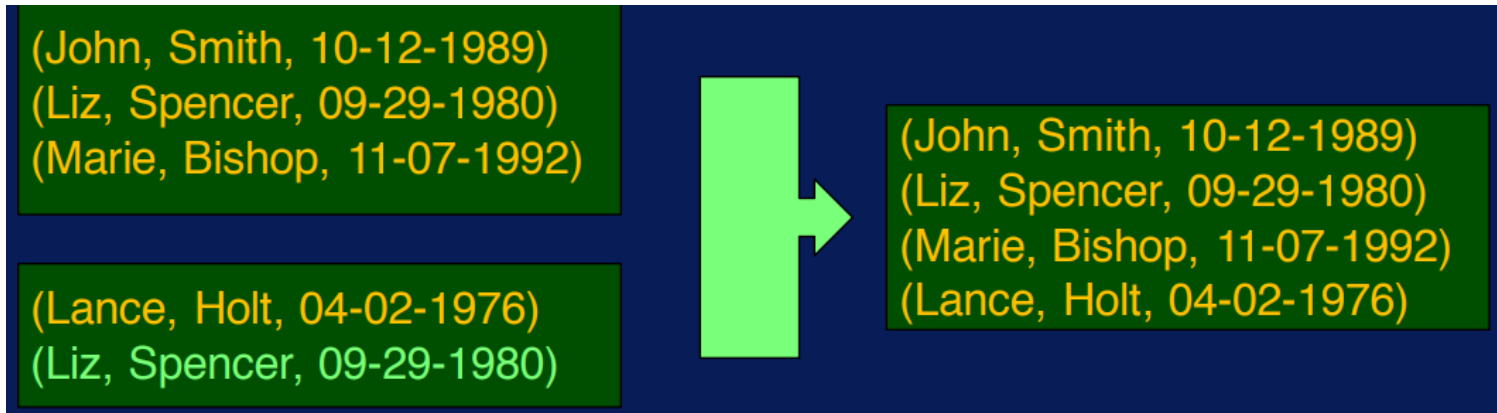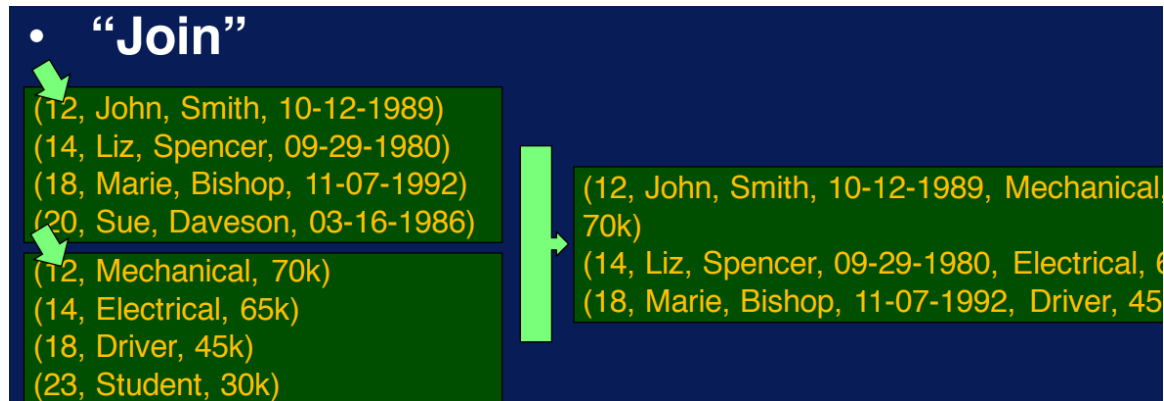| | | |
|---|---|---|
| John, Smith, 10-12-1989, Mechanical, 70000)<br>Liz, Spencer, 09-29-1980, Electrical, 65000)<br>Marie, Bishop, 11-07-1992, Driver, )<br>Steve, Richards, 04-16-1958, 140000) | field 1, field 2 → | (John, Smith)<br>(Liz, Spencer)<br>(Marie, Bishop)<br>(Steve, Richards) |

Operations

- **Union:** given two data sets
  - Create a new data set with elements from the two data sets
  - Erasing duplicates

(John, Smith, 10-12-1989)
(Liz, Spencer, 09-29-1980)
(Marie, Bishop, 11-07-1992)

(Lance, Holt, 04-02-1976)
(Liz, Spencer, 09-29-1980)

(John, Smith, 10-12-1989)
(Liz, Spencer, 09-29-1980)
(Marie, Bishop, 11-07-1992)
(Lance, Holt, 04-02-1976)

**Operations**

- **Join:** given two data sets with complementary structure
    - Create a new group with elements of both data sets
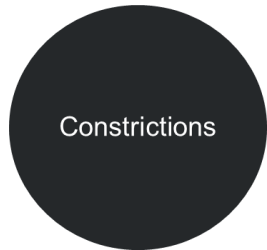    - Erasing the duplicates

- **"Join"**

(12, John, Smith, 10-12-1989)
(14, Liz, Spencer, 09-29-1980)
(18, Marie, Bishop, 11-07-1992)
(20, Sue, Daveson, 03-16-1986)

(12, Mechanical, 70k)
(14, Electrical, 65k)
(18, Driver, 45k)
(23, Student, 30k)

(12, John, Smith, 10-12-1989, Mechanical, 70k)
(14, Liz, Spencer, 09-29-1980, Electrical, 6
(18, Marie, Bishop, 11-07-1992, Driver, 45

Constrictions

- **Constrictions: there are logical propositions thta data must complain. example: each person has only one name**

- **Different models have different ways to express constrictions**

Constrictions

- **Constriction types**

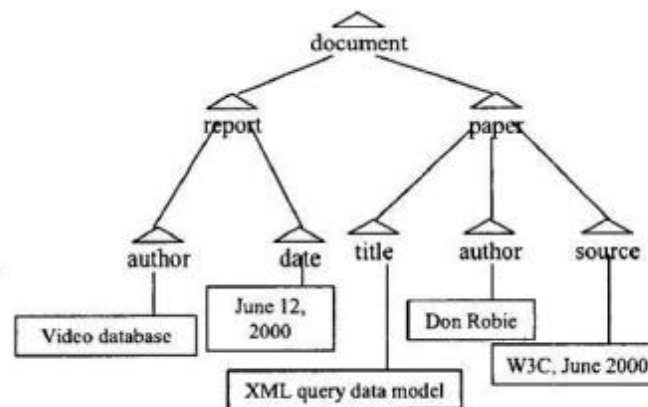| | |
|---|---|
| **Value constrictions:**<br>Age can not be negative | **Unicity:**<br>each person has only one name |
| **Cardinality:**<br>each person has between 1 and 5 pone numbers | **Type:**<br>Surname is Alpha numeric |
| Domain:<br>Days are from 1 to 31 | **Estructural: it sets restrictions to the structure instead to the data**. Example: it must be a2x2 matrix |

Semi structured

- **Semi structured data:**
  - **It is really a tree**
  - The queries actually involve navigating the tree until in common ancestor
  - **Example, query: June 12 2000 and author___document**

```
<document>
  <report>
    <author>Video database</author >
    <date>June 12, 2000</date>
  </report >
  <paper>
    <title>XML query data model</title>
    <author>Don Robie</author>
    <source>W3C, June 2000</source>
  </paper>
</document>
```
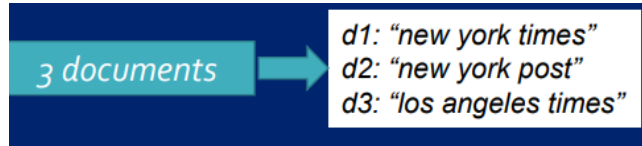
**Unstructured**

- **Vector Space Model: ( we will review deeply in data analysis part next presentations)**
  - **Documents as vectors**

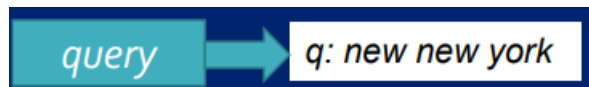| 3 documents | → | d1: "new york times"<br>d2: "new york post"<br>d3: "los angeles times" |
|---|---|---|

|    | angeles | los | new | post | times | york |
|----|---------|-----|-----|------|-------|------|
| d1 | 0 | 0 | 1 | 0 | 1 | 1 |
| d2 | 0 | 0 | 1 | 1 | 0 | 1 |
| d3 | 1 | 1 | 0 | 0 | 1 | 0 |

|    | angeles | los | new | post | times | york | Length |
|----|---------|-----|-----|------|-------|------|--------|
| d1 | 0 | 0 | 0.584 | 0 | 0.584 | 0.584 | 1.011 |
| d2 | 0 | 0 | 0.584 | 1.584 | 0 | 0.584 | 1.786 |
| d3 | 1.584 | 1.584 | 0 | 0 | 0.584 | 0 | 2.316 |

Length of d1 = sqrt(0.584^2+0.584^2+0.584^2)=1.011

- **Queries as vectors for information retrieval**

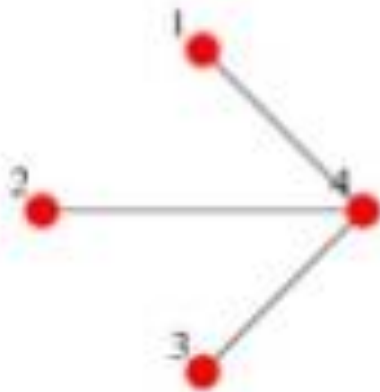| query | → | q: new new york |
|---|---|---|

- **Cosine distance**

$$similarity = \cos(\theta) = \frac{A \cdot B}{\|A\|\|B\|} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}}$$

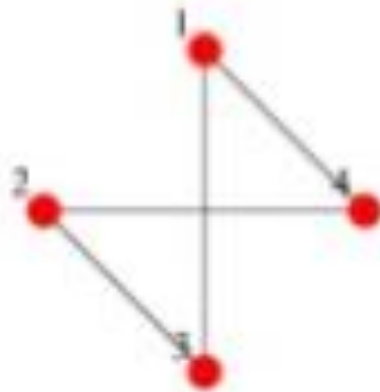cosSim(d1,q) = (0.584*0.584+0.584*0.292) / (1.011*0.652) = 0.776
cosSim(d2,q) = (0.584*0.584) / (1.786*0.652) = 0.292
cosSim(d3,q) = (0.584*0.292) / (2.316*0.652) = 0.112

16

- **Other models: matrices**

- **csv formal does not mean relational model : the format is not the model**

**Book1: Big Data for Dummies**

- Chapter 2:
    - Pg. 25-3

- Chapter 7: Operational data bases pg. 85-100 ( some parts but the relevant content will be review in next presentations)