# Designing Experiments
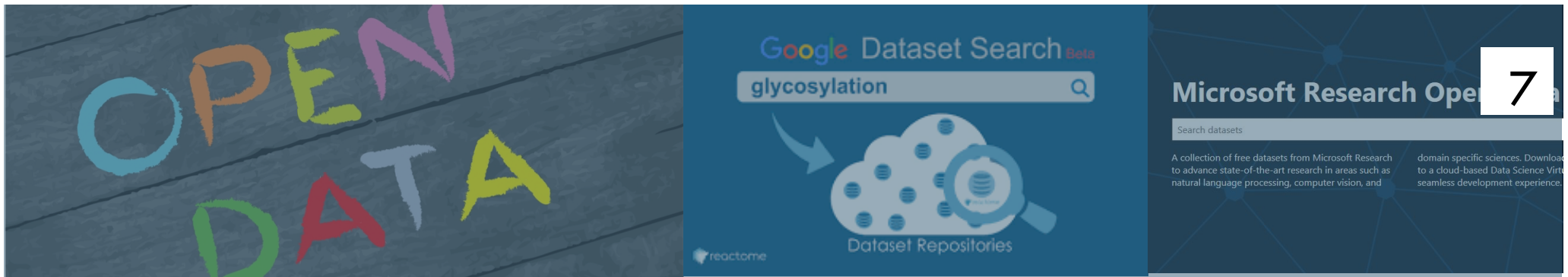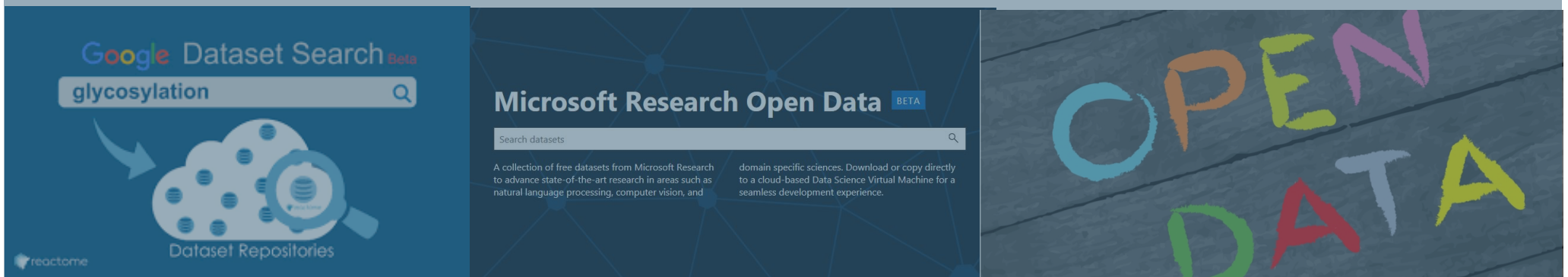## Overview of data management & exploitation solutions

**Genoveva Vargas-Solar**
**French Council of Scientific Research, LIG**
**genoveva.vargas@imag.fr**

http://vargas-solar.com/data-centric-smart-everything/

**Democratized access to open data collections & algorithms**
**released under different conditions & qualities**

# Social Data Science



# Network Science



# Digital humanities



# Computational Science



| Computation (Algorithm: mathematical model) | Experiment setting (Architecture: computing environment) |
|---|---|

**Volume**

**DATA**

**Variety**

**Value**

**Velocity**

**Veracity**

| 1000 Yottabytes | 1 Brontobyte |
|---|---|
| 1000 Brontobytes | 1 Geopbyte |

3

# QUERYING APPROACHES

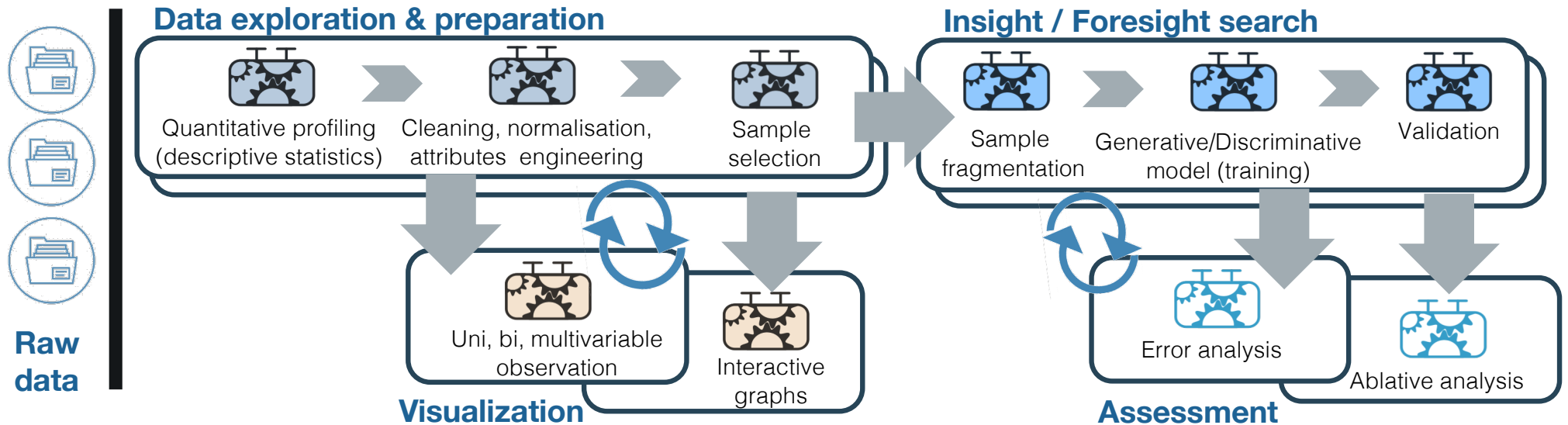| | Databases ¦ information retrieval | | Data Science |
|---|---|---|---|
| Query Types | - Relational, multi-dimensional, spatio-temporal, aggregation<br>- Patterns, regular expressions | - (Dis)conjunctive<br>- Navigational | - Exploratory<br>- Analytics: modelling & predicting |
| Execution model | On demand/continuous | On demand | Step by step |
| Results properties | - Completeness (full/partial)<br>- Fussiness | Approximation<br><br>Precision/recall I Probabilistic | Approximation<br>with some error degree<br>Data, queries, samples |
| Dataset content | Intention model | Extension model | Extension (raw) |
| | *Data structure*<br><br>*table, key-value, tuple, document, graph* | *Quantitative representation*<br><br>*Frequency matrix, Statistical profiling*<br><br>*Semantic representation*<br><br>*Ontology, Terms graph* | *csv, XML, JSON, BLOB, …* |

4

# EXPLORATORY QUERYING

| Data Science |
|---|
| Query Types |
| - Exploratory<br>- Analytics: modelling & predicting |
| Execution model |
| Step by step |
| Results properties |
| Approximation<br>with some error degree<br>Data, queries, samples |
| Dataset content |
| Extension (raw) |

*csv, XML, JSON, BLOB, …*

*Data structure*

*table, key-value, tuple, document, graph*

- Methodologies weaving data management, greedy algorithms

- Programming models that must be tuned to be deployed in different target architectures

Data collections as backbone for conducting **experiments**, drive hypothesis and lead to "valid" conclusions, models, simulations, understanding

# DATA CONSUMPTION PHILOSOPHIES

## Diogenes approach

**Data exploration & preparation**

Quantitative profiling (descriptive statistics)

Cleaning, normalisation, attributes engineering

Sample selection

**Insight / Foresight search**

Sample fragmentation

Generative/Discriminative model (training)

Validation

**Visualization**

Uni, bi, multivariable observation

Interactive graphs

**Assessment**

Error analysis

Ablative analysis

**Raw data**

- **Data**: observations of phenomena often described as series of features/attributes
- **Query**: analytics objective (looks for insights or foresights) expressed as a pipeline of operations guided by the conditions and characteristics of the data
- **Result**: a model or prediction with associated assessment indexes, not definitive accepted with an associated error margin, accepted by comparison

# DATA PROCESSING AND ANALYSIS

"Does it make sense to invest in low-carbon technologies?"



IoT

**SPATIO – TEMPORAL SERIES**

**Inhouse observed variables**

- Electric consumption
- Indoor temperature
- Indoor humidity
- Gas consumption
- Outdoor temperature
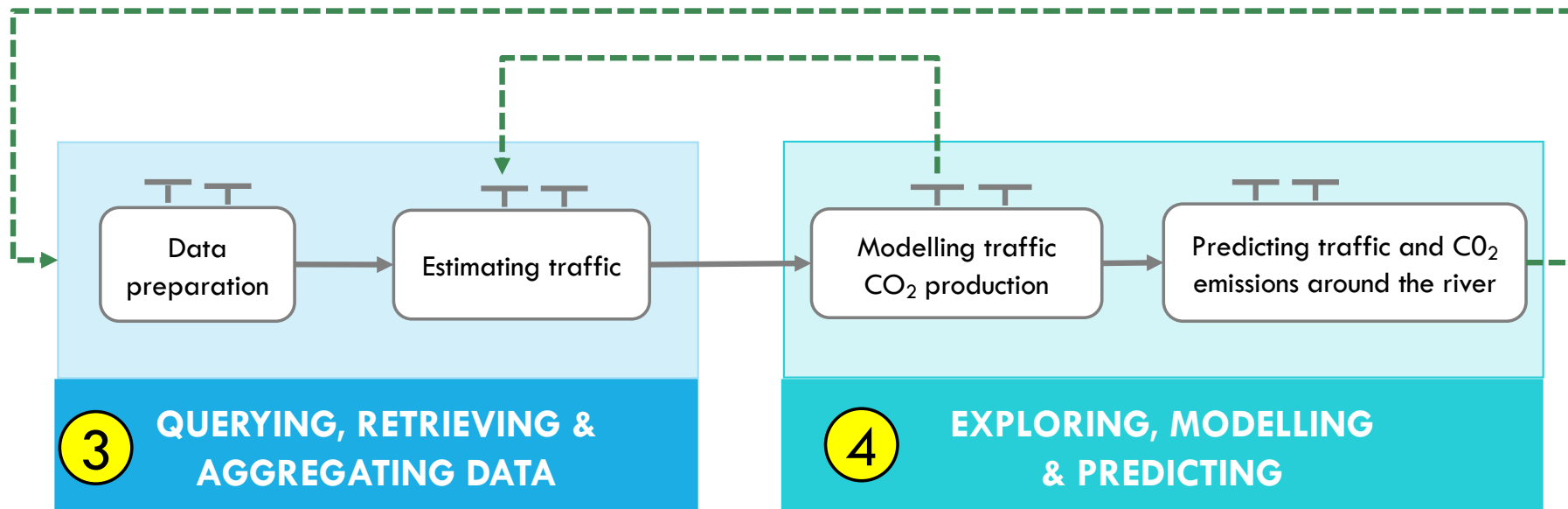- Outdoor humidity

**Meteorological variables**

- Total precipitation
- Total cloud cover
- Shortwave radiation
- Wind direction
- Snowfall amount
- Sunshine duration
- Wind speed

- Heterogeneous masses of data , often spatio-temporal and produced as streams (*i.e.*, spatio-temporal series),
- Backbone of analytical and prediction processes
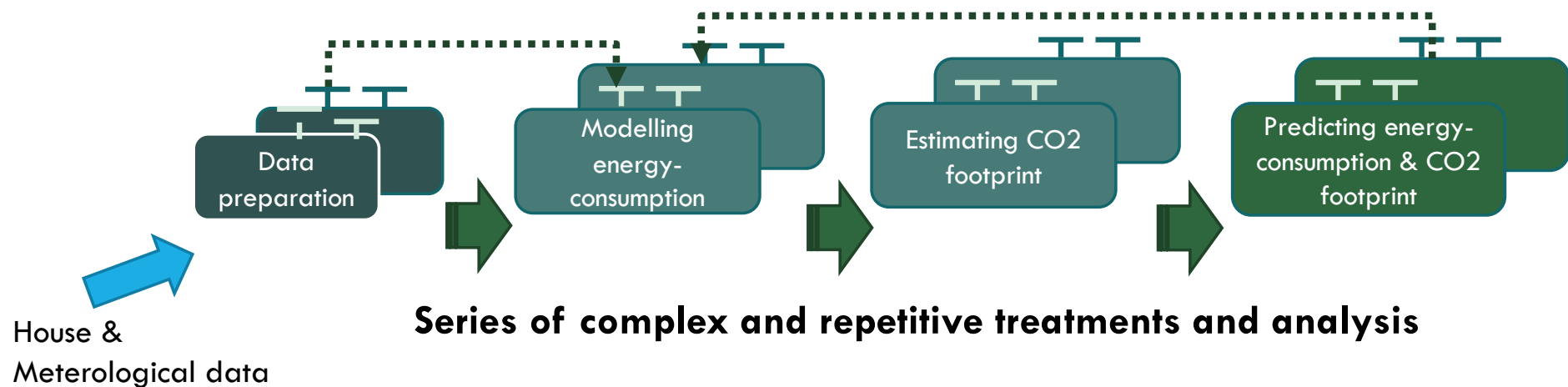
# DATA SCIENCE PIPELINE

(1) COMPLEX AND REPETITIVE PROCESSING & ANALYSIS TASKS

**Example:** (2) *Will jogging in the morning around the river next weekend reduce breathing $CO_2$ ?*
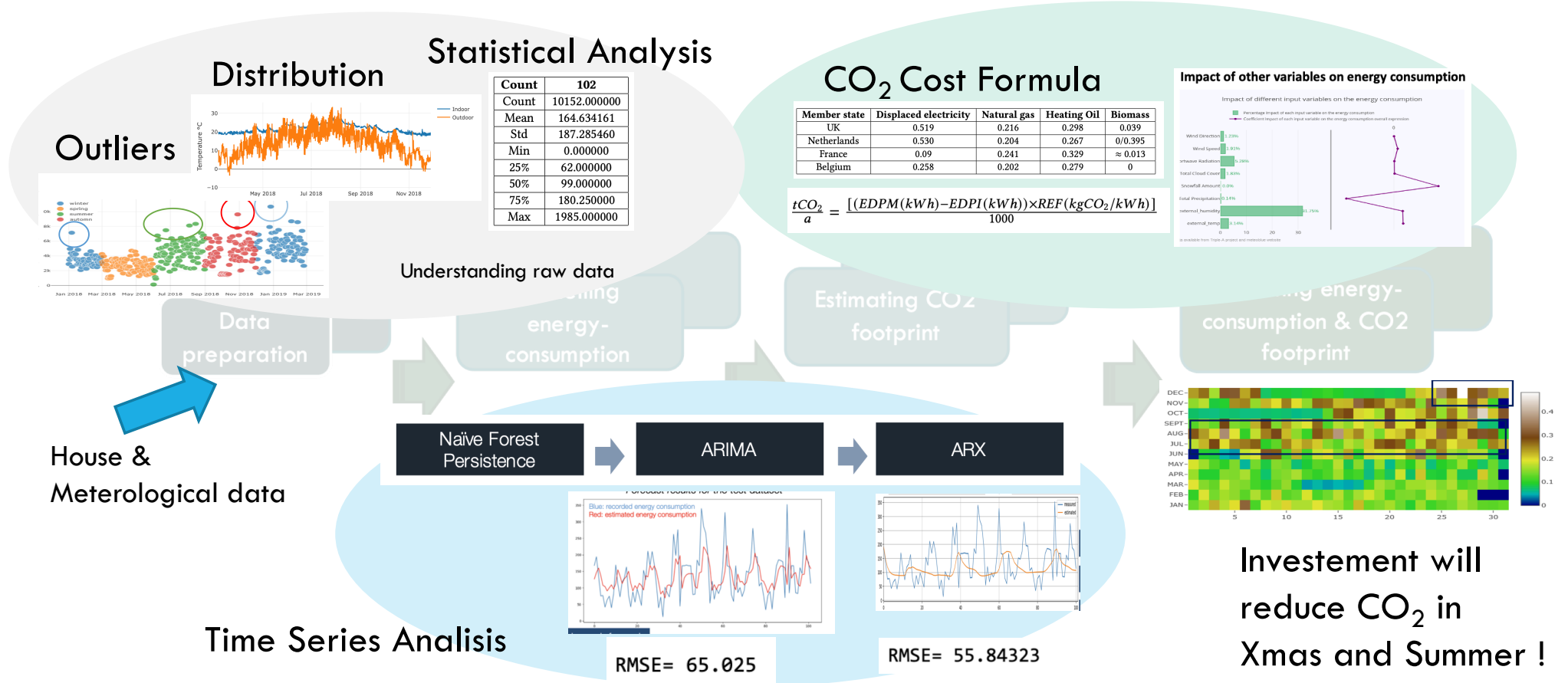
Data preparation → Estimating traffic → Modelling traffic $CO_2$ production → Predicting traffic and $CO_2$ emissions around the river

(3) QUERYING, RETRIEVING & AGGREGATING DATA

(4) EXPLORING, MODELLING & PREDICTING

# DATA SCIENCE PIPELINE

"Does it make sense to invest in low-carbon technologies?"

Data preparation

Modelling energy-consumption

Estimating CO2 footprint

Predicting energy-consumption & CO2 footprint

House & Meterological data

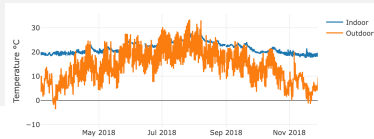**Series of complex and repetitive treatments and analysis**

- **Exploration & preparation:** Identify which variables impact my household $CO_2$ footprint
- **Analysis:** Predict whether investing in low-carbon technologies will decrease my $CO_2$ footprint
- **Assessment:** Identify which prediction model is better for my household energy consumption pattern
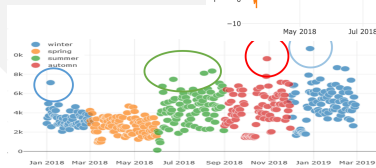
# DATA SCIENCE PIPELINE

## Statistical Analysis

### Distribution

### Outliers

| Count | 102 |
|---|---|
| Count | 10152.000000 |
| Mean | 164.634161 |
| Std | 187.285460 |
| Min | 0.000000 |
| 25% | 62.000000 |
| 50% | 99.000000 |
| 75% | 180.250000 |
| Max | 1985.000000 |

Understanding raw data

## CO₂ Cost Formula

**Impact of other variables on energy consumption**

| Member state | Displaced electricity | Natural gas | Heating Oil | Biomass |
|---|---|---|---|---|
| UK | 0.519 | 0.216 | 0.298 | 0.039 |
| Netherlands | 0.530 | 0.204 | 0.267 | 0/0.395 |
| France | 0.09 | 0.241 | 0.329 | ≈ 0.013 |
| Belgium | 0.258 | 0.202 | 0.279 | 0 |

$$\frac{tCO_2}{a} = \frac{[(EDPM(kWh) - EDPI(kWh)) \times REF(kgCO_2/kWh)]}{1000}$$

Data preparation

Modelling energy-consumption

Estimating CO2 footprint

Modelling energy-consumption & CO2 footprint

House & Meterological data

Naïve Forest Persistence → ARIMA → ARX

**Time Series Analisis**

Blue: recorded energy consumption
Red: estimated energy consumption

RMSE= 65.025

RMSE= 55.84323

Investement will reduce CO₂ in Xmas and Summer !

# DATA SCIENCE PIPELINE

**(1)** **Artisanal design**
depending on
data scientist/engineers
*"expertise"*

**(2)** In-house programming using many
different **libraries, stacks, tools
difficult to integrate**



**DATA SCIENCE LABS**

- kaggle.com
- Google Colab
- Azure Notebooks

**DATA SCIENCE STACKS**

**BD SERVICES PLATFORMS**

# DATA SCIENCE PIPELINE

## CHALLENGES

**1** Efficient execution on distributed architectures requiring **important engineering effort**

**2** **Tuning** for improving **data management** across nodes hosting software components and libraries

### DATA SCIENCE LABS

kaggle.com

Google Colab

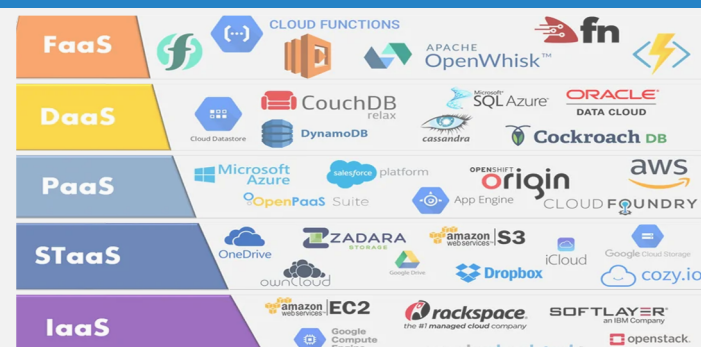Azure Notebooks

### DATA SCIENCE STACKS



### BD SERVICES PLATFORMS

# Machine Learning Studios

## Pipeline ①

**Task:**
**Function**
**I/O constraints** ②

College - Copy.csv
Select Columns in Dataset ✓
Linear Regression
Split Data ✓
Train Model ✓
Score Model ✓
Evaluate Model

### ③ Tracking
Record & query experiments: code, configs, results .. Etc

### ④ Projects
Packaging format for reproducible runs on any platform

### ⑤ Models
General model format that supports diverse deployment tools

## DS Tools

Microsoft Azure Machine Lear

Search experiment items 🔍
- Trained Models
- Data Format Conversions
- Data Input and Output
- Data Transformation
- Feature Selection
- Machine Learning
- OpenCV Library Modules
- Python Language Modules
- R Language Modules
- Statistical Functions
- Text Analytics
- Time Series

| Services | | | Tools |
|---|---|---|---|
| **Custom AI** | **Prebuilt AI** | **Conversational AI** | **Coding and management tools** |
| Azure ML | Cognitive Services | Bot Framework | VS AI Tools / Azure ML Studio / Azure ML Workbench |

Others (Pycham and Jupyter Notebooks …)

Deep learning frameworks

| **Infrastructure** | | | | | | | Cognitive Toolkit | Tensorflow | Caffe |
|---|---|---|---|---|---|---|---|---|---|
| AI on data | | | | AI compute | | | | | |
| Cosmos DB | SQL DB | SQL DW | Data Lake | Spark | DSVM | Batch AI | ACS | | |
| CPU, FPGA, GPU | | | | | | | | | |

Others (Scikit-Learn, MxNet, Keras, Chainer, Gluon,…)

**Enactment (AI Platform, e.g. Microsoft, Databricks MLFlow, Google AI)**

13

# Cloud ML Studios

| | Amazon | Microsoft | Google | IBM |
|---|---|---|---|---|
| **Automated and semi-automated ML services** | | | | |
| | AmazonML | MS AzureML Studio | Cloud AutoML | IBM Watson Model Builder |
| Classification | ✓ | ✓ | ✓ | ✓ |
| Regression | ✓ | ✓ | ✓ | ✓ |
| Clustering | ✓ | ✓ | ✗ | ✗ |
| Anomaly detection | ✗ | ✓ | ✗ | ✗ |
| Recommendation | ✗ | ✓ | ✓ | ✗ |
| Ranking | ✗ | ✓ | ✗ | ✗ |
| **Platforms for custom modelling** | | | | |
| | Amazon SageMaker | Azure ML Services | Google ML Engine | IBM WatsonM Studio |
| Built-in algorithms | ✓ | ✗ | ✓ | ✓ |
| Supported Frameworks | Tensorflow, MXNet, Keras, Gluon, Pytorch, Caffe2, Chainer, Torch | Tensorflow, Scikit-Learn, MS Cognitive Toolkit, SparkML | Tensorflow, Scikit-Learn, XGBoost, Keras | Tensorfoow, SparkMLib, Scikit-Learn, XGBoost, PyTorch, IBM SPSS, PMML |

# PREPARING DATA

# DATA MODELS

## Tuple

- Row in a relational table, where attributes are pre-defined in a schema, and the values are scalar

## Document

- Allows values to be nested documents or lists, as well as scalar values.

- Attributes are not defined in a global schema

## Extensible record

- Hybrid between tuple and document, where families of attributes are defined in a schema, but new attributes can be added on a per-record basis

# DATA STORES

Key-value

- Systems that store values and an index to find them, based on a key

Document

- Systems that store documents, providing index and simple query mechanisms

Extensible record

- Systems that store extensible records that can be partitioned vertically and horizontally across nodes

Graph

- Systems that store model data as graphs where nodes can represent content modelled as document or key-value structures and arcs represent a relation between the data modelled by the node

Relational

- Systems that store, index and query tuples

# KEY STORE VALUES

"Simplest data stores" use a data model similar to the memcached distributed in-memory cache

Single key-value index for all data

Provide a persistence mechanism

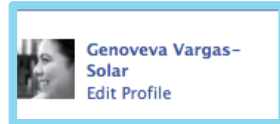Replication, versioning, locking, transactions, sorting

API: inserts, deletes, index lookups

No secondary indices or keys

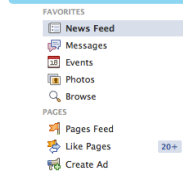| System | Address |
| --- | --- |
| Redis | code.google.com/p/redis |
| Scalaris | code.google.com/p/scalaris |
| Tokyo | tokyocabinet.sourceforge.net |
| Voldemort | project-voldemort.com |
| Riak | riak.basho.com |
| Membrain | schoonerinfotech.com/products |
| Membase | membase.com |

```
SELECT    name, pic, profile_url
FROM      user
WHERE     uid = me()
```
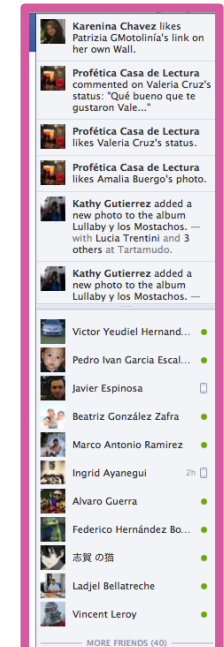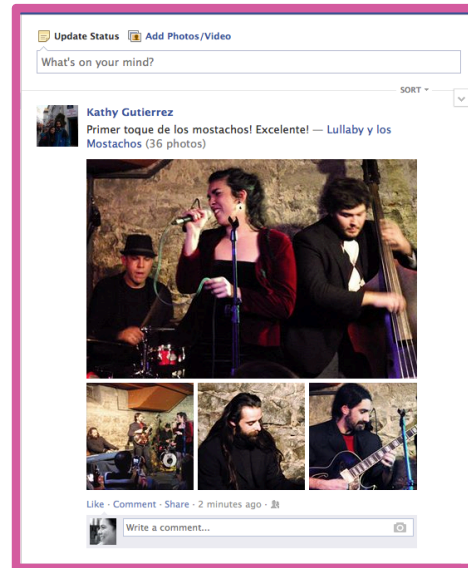
Genoveva Vargas-Solar
Edit Profile

```
SELECT    message, attachment
FROM      stream
WHERE     source_id = me() AND type = 80
```

Karenina Chavez likes
Patrizia GMotolinía's link on
her own Wall.

Profética Casa de Lectura
commented on Valeria Cruz's
status: "Qué bueno que te
gustaron Vale..."

Profética Casa de Lectura
likes Valeria Cruz's status.

Profética Casa de Lectura
likes Amalia Buergo's photo.

Kathy Gutierrez added a
new photo to the album
Lullaby y los Mostachos. —
with Lucia Trentini and 3
others at Tartamudo.

Kathy Gutierrez added a
new photo to the album
Lullaby y los Mostachos. —

Victor Yeudiel Hernand...
Pedro Ivan Garcia Escal...
Javier Espinosa
Beatriz González Zafra
Marco Antonio Ramirez
Ingrid Ayanegui          2h
Alvaro Guerra
Federico Hernández Bo...
志賀 の猫
Ladjel Bellatreche
Vincent Leroy
MORE FRIENDS (40)

FRIENDS
  Close Friends
  Family
  National Laboratory on ...
  UDLA, Universidad de la...
  Colegio Humboldt Puebla
  Fundación Universidad d...
  Grenoble, France Area
  Colleagues

FAVORITES
  News Feed
  Messages
  Events
  Photos
  Browse
PAGES
  Pages Feed
  Like Pages        20+
  Create Ad

```
SELECT    name
FROM      friendlist
WHERE     owner = me()
```

Update Status    Add Photos/Video
What's on your mind?
                                SORT

Kathy Gutierrez
Primer toque de los mostachos! Excelente! — Lullaby y los
Mostachos (36 photos)

Like · Comment · Share · 2 minutes ago ·
Write a comment...

GROUPS
  Egresados UDLAP
  Such Good People – an i...
  Monis – groupe de soutien
  Découvre ce film qui s'e...
  AIDONS LE REFUGE
  Create Group...

```
SELECT    name
FROM      group
WHERE     gid IN ( SELECT  gid
                   FROM    group_member
                   WHERE   uid = me() )
```
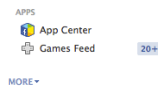
```
SELECT    name, pic
FROM      user
WHERE     online_presence = "active"
AND
          uid IN ( SELECT  uid2
                   FROM    friend
                   WHERE   uid1 = me() )
```

APPS
  App Center
  Games Feed        20+
MORE

https://developers.facebook.com/docs/reference/fql/

19

# DOCUMENT STORES

Support more complex data: pointerless objects, i.e., documents

Secondary indexes, multiple types of documents (objects) per database, nested documents and lists, e.g. B-trees

Automatic sharding (scale writes), no explicit locks, weaker concurrency (eventual for scaling reads) and atomicity properties

API: `select, delete, getAttributes, putAttributes` on documents

Queries can be distributed in parallel over multiple nodes using a map-reduce mechanism

| SYSTEM | ADDRESS |
|--------|---------|
| SimpleDB | amazon.com/simpledb |
| Couch DB | couchdb.apache.org |
| Mongo DB | mongodb.org |
| Terrastore | code.google.com/terrastore |

# DOCUMENT STORES

```
{
    "name": "Genoveva Vargas-Solar",
    "id": "805114856"
}
```

```
{
"data": [
  {
    "name": "Genoveva Vargas-Solar",
    "pic": "https://fbcdn-profile-a.akamaihd.net/hprofile-ak-ash4/275915_805114856_16986061_s.jpg",
    "profile_url": "https://www.facebook.com/genoveva.vargas"
  }
]
}
```

```
{
"data": [
  {
    "name": "$$$ Se Vende Jeep Compass 2008 - 60,000kms. $$$"
  },
  {
    "name": "Découvre ce film qui s'engage pour le mariage pour tous"
  },
  {
    "name": "emepink"
  },
  {
    "name": "Such Good People - an indie screwball comedy"
  },
  {
    "name": "Comunidad Mexicana de Tecnologías Semánticas"
  },
  {
    "name": "TI-502 Administración de Datos"
  },
  {
    "name": "exaUDLAP Sistemas Computacionales"
  },
  {
    "name": "\"Hombre Nuevo\" artículos de valores humanos del P. Otaolaurruchi"
  },
  {
    "name": "LACCIR"
  },
  {
    "name": "Monis - groupe de soutien"
  },
  {
    "name": "Red Temática de las TIC"
  },
```

```
{
"data": [
  {
    "message": "",
    "attachment": {
      "media": [
        {
          "href": "https://www.facebook.com/photo.php?fbid=10151871935952502&set=a.99396912501.109184.98871212501&type=1",
          "alt": "",
          "type": "photo",
          "src": "https://fbcdn-photos-e-a.akamaihd.net/hphotos-ak-ash3/1146527_10151871935952502_258686255_s.jpg",
          "photo": {
            "aid": "98871212501_109184",
            "pid": "98871212501_1073742168",
            "fbid": "10151871935952502",
            "owner": 98871212501,
            "index": 1,
            "width": 611,
            "height": 458,
            "images": [
              {
                "src": "https://fbcdn-photos-e-a.akamaihd.net/hphotos-ak-ash3/1146527_10151871935952502_258686255_s.jpg",
                "width": 130,
                "height": 97
              }
            ]
          }
        }
      ],
      "name": "Timeline Photos",
      "href": "https://www.facebook.com/album.php?fbid=99396912501&id=98871212501&aid=109184",
      "caption": "El sutil arte de cantinflear.\r\n\r\nvía - Lectura Cinematográfica",
      "description": "",
```

# EXTENSIBLE RECORD STORES

Basic data model is rows and columns

Basic scalability model is splitting rows and columns over multiple nodes

- Rows split across nodes through sharding on the primary key
  - Split by range rather than hash function
  - Rows analogous to documents: variable number of attributes, attribute names must be unique
  - Grouped into collections (tables)
  - Queries on ranges of values do not go to every node

Columns are distributed over multiple nodes using "column groups"

- Which columns are best stored together
- Column groups must be pre-defined with the extensible record stores

| SYSTEM | ADDRESS |
|--------|---------|
| HBase | hbase.apache.com |
| HyperTable | hypertable.org |
| Cassandra | incubator.apache.org/cassandra |

# SCALABLE RELATIONAL SYSTEMS

SQL: rich declarative query language

Databases reinforce referential integrity

ACID semantics

Well understood operations:

- Configuration, Care and feeding, Backups, Tuning, Failure and recovery, Performance characteristics

Use small-scope operations

- Challenge: joins that do not scale with sharding

Use small-scope transactions

- ACID transactions inefficient with communication and 2PC overhead

Shared nothing architecture for scalability

Avoid cross-node operations

| System | Address |
|--------|---------|
| MySQL C | mysql.com/cluster |
| Volt DB | voltdb.com |
| Clustrix | clustrix.com |
| ScaleDB | scaledb.com |
| Scale Base | scalebase.com |
| Nimbus DB | nimbusdb.com |

# 1970 - 2000 RELATIONAL DB

*More than 30 years: maturity!*

Theoretical & Practical aspects (DBMS)

Domains & $R \subseteq D_1 \times D_2 \times \ldots D_n$, Algebra $\rightarrow$

1st Order Predicate Logic

Languages: SQL (wins), QUEL, QBE

DBMS Prototypes (1975), Products (1980)

A major improvement in DB: provide **data independence** & a simple, **tabular view** of data

Normal Forms & Dependencies (DB design, **consistency**)

Controversial: missing values, duplicates
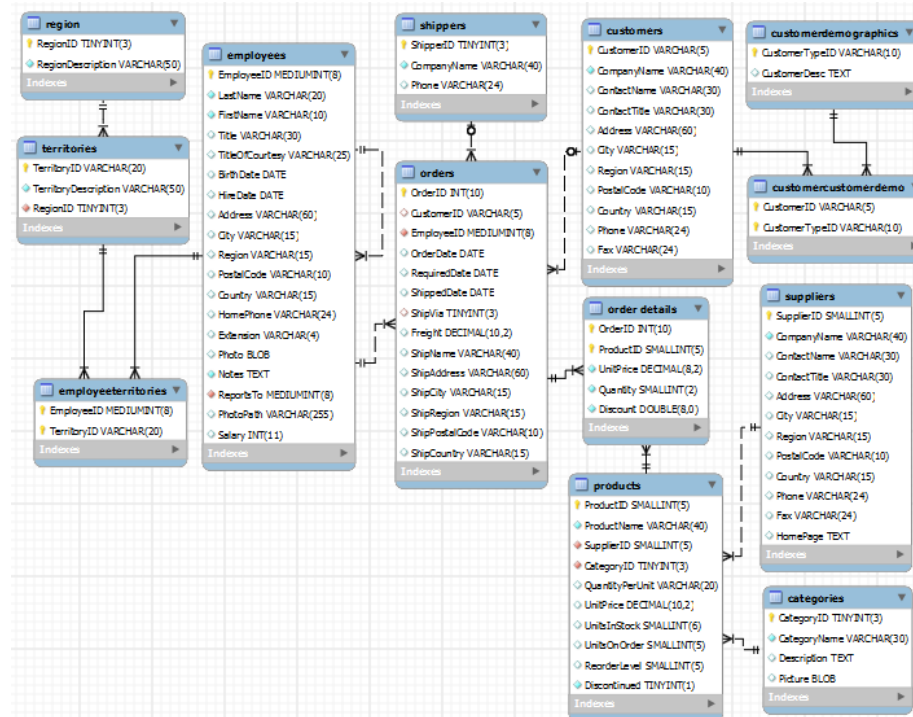
$R \times S$

$R \cup S$

$R - S$

$R[\alpha]$

$R : \varphi$

-------

$R * S$

# MODELING DATA COLLECTIONS

Raw data collections        tabular (csv, excel)              - **Relational**

{
  "geometry": {
    "type": "Point",
    "coordinates": [
      4.821773,
      45.7513
    ]
  },
  "full_location": "Autoroute du Soleil,
69005 Lyon",
  "_id": "Criter11185353",
  "properties": {
    "confidentiality": "noRestriction",
    "probability": "certain",
    "mobility": "",
    "creationtime": "2016-06-07 19:40:00",
    "publiceventtype": "",
    "networkmanagementtype": "",
    "observationtime": "2016-06-07
19:40:00",
    "last_update": "2016-06-07 19:43:30",
    "numberoflanesrestricted": "0",
    "effectonroadlayout": "",
    "creator": "CRITER",
    "id": "Criter11185353",

    "firstsupplierversiontime": "2016-06-07
19:40:00",
    "version": "1",
    "linkname": "",
    "type": "VehicleObstruction",
    "status": "active",
    "direction": "bothWays",
    "locationtype": "nonLinkedPoint",
    "disturbanceactivitytype": "",
    "last_update_fme": "2016-06-07
19:44:29",
    "endtime": "",
    "creationreference": "",
    "informationstatus": "real",
    "townname": "Voie Rapide Urbaine de
Lyon",
    "publiccomment": "Bouchon, km 455|Voie
Rapide Urbaine de Lyon",
    "roadmaintenancetype": "",
    "versiontime": "2016-06-07 19:43:26",
    "starttime": "2016-06-07 19:40:00",
    "gid": "39258",
    "abnormaltraffictype": ""
  }
}

- **Key-Value**

Media (XML, JSON, BLOB)     - Column oriented Tabular
            Graph
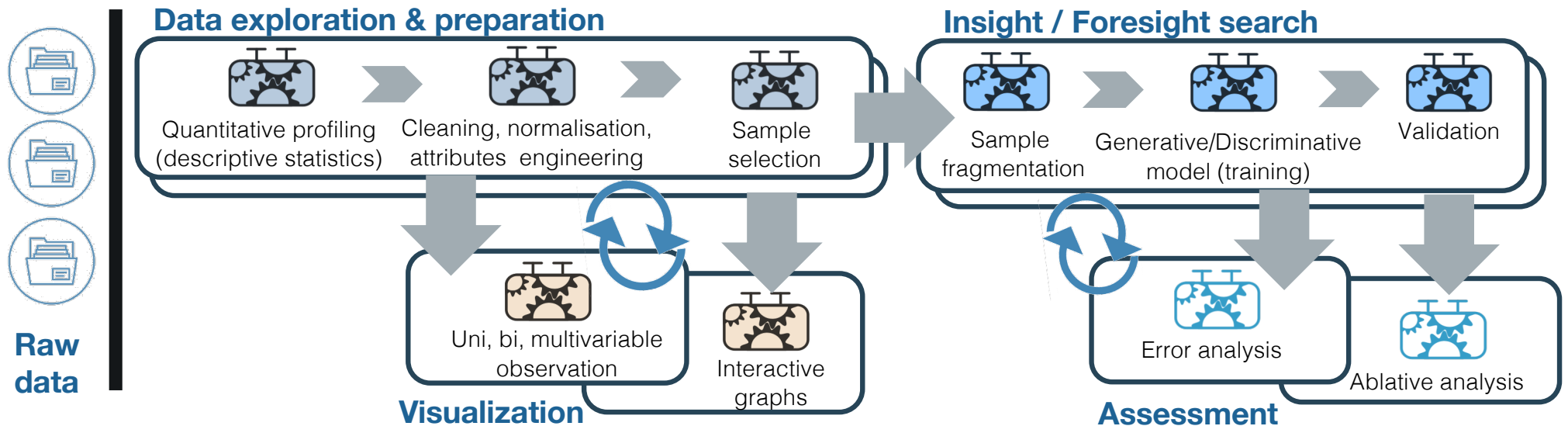
- **Document oriented**

- How to transform data collections ?
- Which is the best adapted model?

→ **Polyglot persistence**

**Approaches dealing with transformation rules inspired in the relational case**

# EXPERIMENT DESIGN

Diogenes approach

**Data exploration & preparation**

Quantitative profiling (descriptive statistics)

Cleaning, normalisation, attributes engineering

Sample selection

**Insight / Foresight search**

Sample fragmentation

Generative/Discriminative model (training)

Validation

**Raw data**

Uni, bi, multivariable observation

Interactive graphs

**Visualization**

Error analysis

Ablative analysis

**Assessment**

# PREPARING DATA

**Obtaining the data**: Read from a file or obtained by scraping the web

**Parsing the data**: Format the data which can be in plain text, fixed columns, CSV, XML, HTML, etc.

**Cleaning the data**: A simple strategy is to remove or ignore incomplete records

**Building data structures**: A data structure that lends itself to the analysis we are interested in.

Databases provide a mapping from keys to values, so they serve as dictionaries

# ANALYSING INCOME ACCORDING TO GENDER

Financial parameters related to the US population*

- **Features**: Age, sex, marital, country, income, education, occupation, capital gain, etc.

- **Question**: Are men more likely to become high-income professionals than women, i.e., to receive an income of over $50,000 per year?

- Preparing data collections

  - Read and check the data

  - Represent the data, for instance using a tabular data structure with features (columns) and records (rows)

  - Group the data

* UCI's Machine Learning Repository: https://archive.ics.uci.edu/ml/datasets/Adult

# EXPLORATORY DATA ANALYSIS

Measurements and categories represent a sample distribution of a variable:

- which approximately represents the population distribution of the variable

- to make tentative assumptions about the population distribution

Different *techniques*:

- **Summarizing the data**

- **Data distributions**

- **Outlier treatment**

- **Kernel density**

https://www.kaggle.com/robikscube/hourly-energy-consumption

SMART ENERGY

# UNDERSTANDING ENERGY CONSUMPTION IN THE PHILIPPINES

https://www.kaggle.com/ljvmiranda/philippines-energy-use

1. What percentage of the population has access to electricty?
   - Access to electricity over time
   - Comparison to South-East Asian (SEA) countries
2. What constitutes my country's energy mix?
   - Energy Mix in the Philippines
   - Comparison to South-East Asian (SEA) countries
     - Fossil-Fuel use
     - Renewable Energy Adoption
3. How are we consuming our energy?
   - Electric Power Consumption over time
   - Consumption footprint

HANDS ON