# Data Mining
# Algorithms and Applications

**Eduardo Ogasawara**

eogasawara@cefet-rj.br
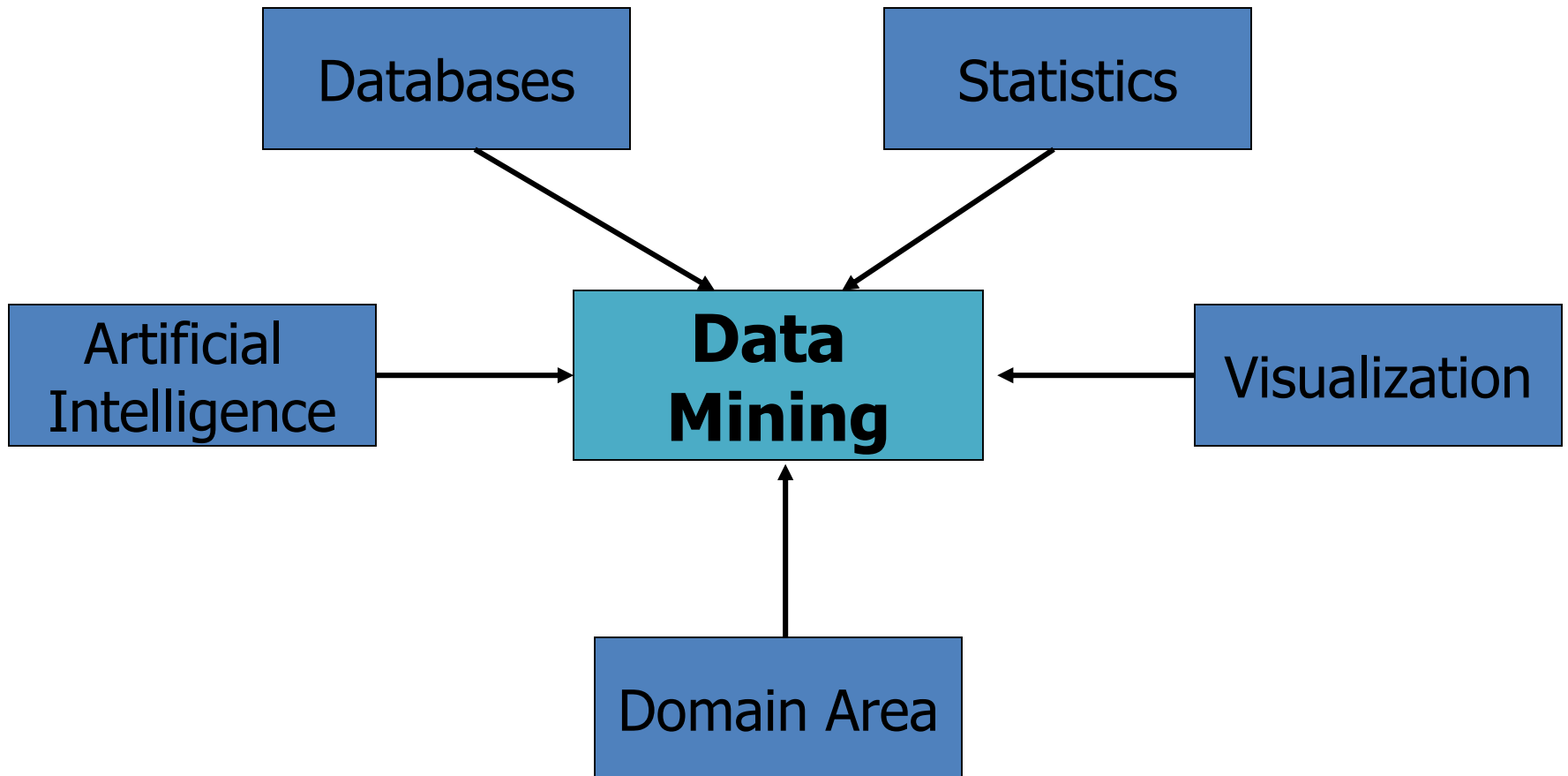
**Centro Federal de Educação Tecnológica**
**Celso Suckow da Fonseca**
**CEFET/RJ**

# *Data Mining*

- **Extraction knowledge (rules, patterns, constraints) from "large" databases**

- **Data Deluge & Big Data introduce more complexity**
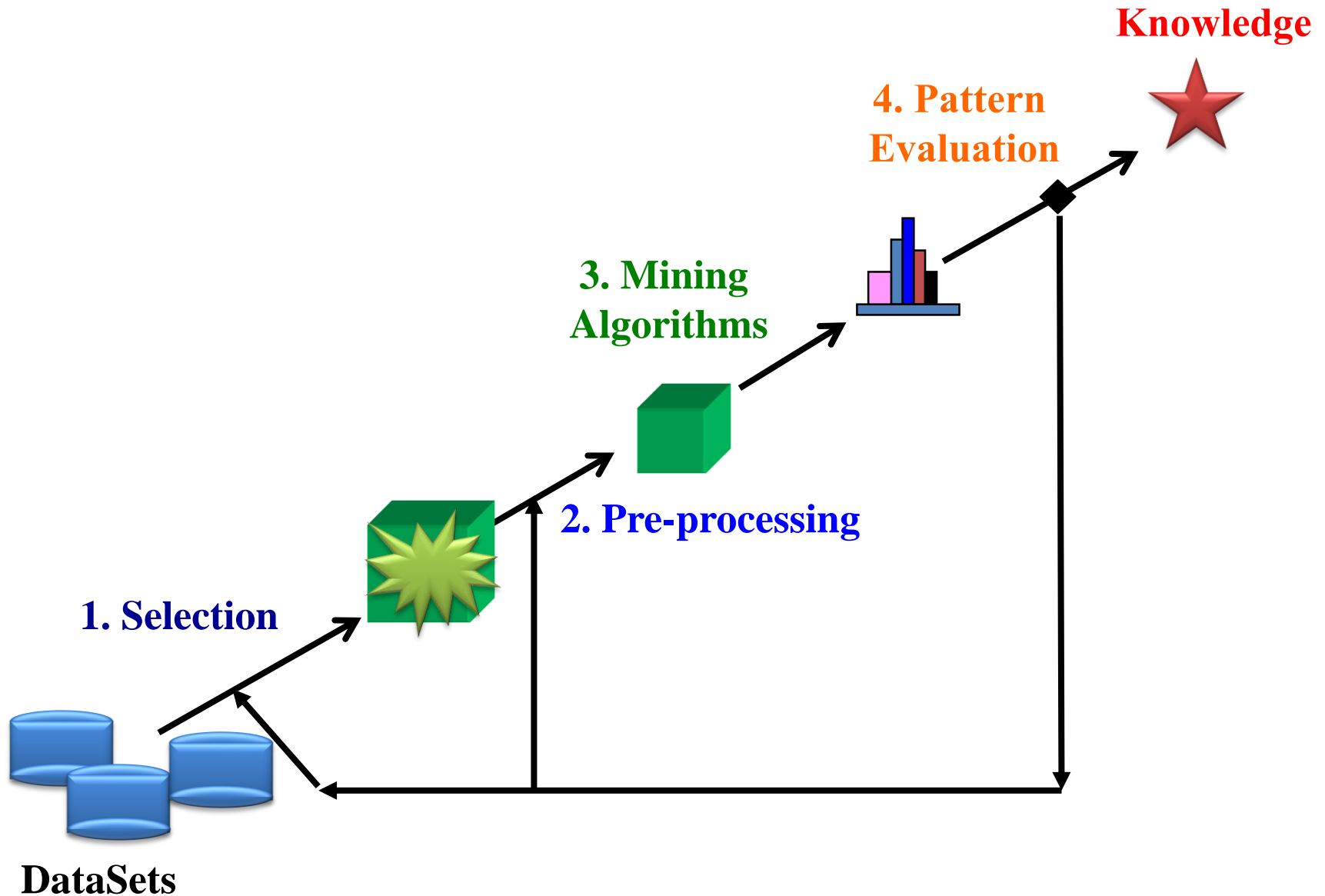  - **We are immersed in data, but hungry for knowledge!**

# *Goal: Answer key questions*

- **Consider an example of a internet sales system**

- **Do we have the data?**
  - **Transactions involving credit cards, loyalty cards, central customer service, lifestyle customer**

- **Can we identify groups?**
  - **Identified clusters for "models" of customers who share certain characteristics: interest, income, spending habits**

- **Can we classify or predict?**
  - **Identified patterns of consumption?**
  - **Can establish a pattern of consumption for single and married?**

- **Can we associate patterns?**
  - **Association / correlation between product sales**
  - **Prediction based on association information**

# *Data Mining Process (DMP)*

**Knowledge**

**4. Pattern Evaluation**

**3. Mining Algorithms**

**2. Pre-processing**

**1. Selection**

**DataSets**

- **Goal**
  - **Develop Data Mining Algorithms and Application for practical problems**

- **Research Process**
  - **Basic Research**
    - **Produce data knowledge extraction algorithms or tools without immediate application for real problems**
  - **Applied Research**
    - **Usage of know or developed algorithms or tools for solving real problems**

| Basic Research | Applied Research |
| --- | --- |

# *Basic Research*

- **Data Management & Data Mining Process**
  - **Data Streaming**
  - **Data Sources Location**
  - **Uniform Data Model**
  - **Data Experiment Representation**
  - **Modeled through Workflows**
- **Algorithms and Methods**
  - **Research on specific Data Mining Activities**

# *DMP: Step #1 – Selection and Integration*

- **Types of data**
  - **Relational, Spatial, Temporal, Textual, Graph**
- **Sources**
  - **Centralized**
  - **Distributed**
  - **Streaming**
- **Integration**
  - **Data integration from different sources**
- **Dataset Production**
    - **Projection**
    - **Selection**

- **Very important step**
  - **60% of Data Mining effort!**
- **Data Reduction**
  - **Filtering or Sampling**
- **Attribution Reduction**
  - **Selection of important attributed**
  - **Reduction of Dimensionality**
- **Transformation**
  - **Combination**
  - **Derivation**
  - **Aggregation**

# *DMP: Step #3 – Data Mining Algorithms and Methods*

- **Focus: Searching for patterns of interest**
- **Data Mining Methods**
  - **Clustering**
  - **Classification and Prediction**
  - **Association**
  - **New Approaches**

# *DMP: Step #4 – Pattern Evaluation*

- **Error Measurements**
  - **Error Metrics**

- **Evaluation of patterns and extracted knowledge**
  - **Visualization**
  - **Transformation**
  - **Redundant pattern removal**

- **Use of discovered knowledge**

- **Are all discovered patterns important?**

# *Applied Research*

- **Focus on solving target application areas (datasets)**
- **Basic domain of target area**
  - **data**
  - **Problem**
  - **Identification of goals**
- **Examples**
  - **Astronomy**
  - **Production engineering**
  - **Social Networks**
  - **Computing in Educational**

(GC)　　　(GR)

Forecasts from ARIMA(3,2,4)

# *Forecast Sea Surface Temperature in Atlantic Pacific Ocean*

**NGC 7793: a spiral galaxy**

# *Prediction of electrical power output of a base load operated combined cycle power plan*



**Tüfekci, P. 2014**

# *Referências*