# Noise and Outlier Detection

BORUT SLUBAN

*DATA MINING AND KNOWLEDGE DISCOVERY*

# Anomalies?

▶ Errors in the data – noise

  ▶ Animals of white color



▶ Exceptions or Outliers

  ▶ Herd of sheep

# Motivation

- **Noise** in data negatively affect data mining results.

  (Zhu et al., 2004)

- False medical diagnosis (**classification noise**) can have serious consequences

  (Gamberger et al. 2003)

- **Outlier** detection proved to be effective in detection of network intrusion and bank fraud.

  (Aggarwal and Yu, 2001)

# Detecting noise and outliers

- Used for:

  - Improving machine learning performance through cleaning of training data

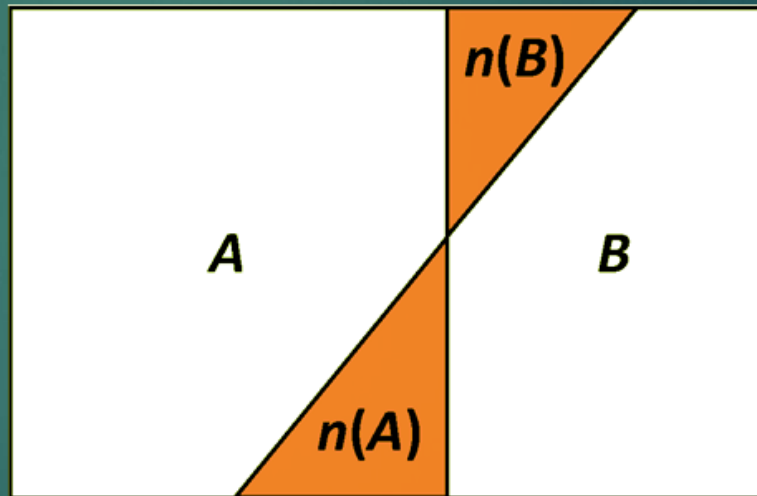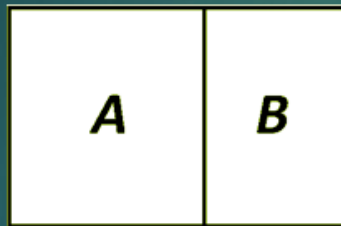  - Data understanding and knowledge expansion by discovering potentially interesting exceptional cases in data

# Detecting noise and outliers

- Nature

    - Follows certain patters

    - Adheres to the laws of physics
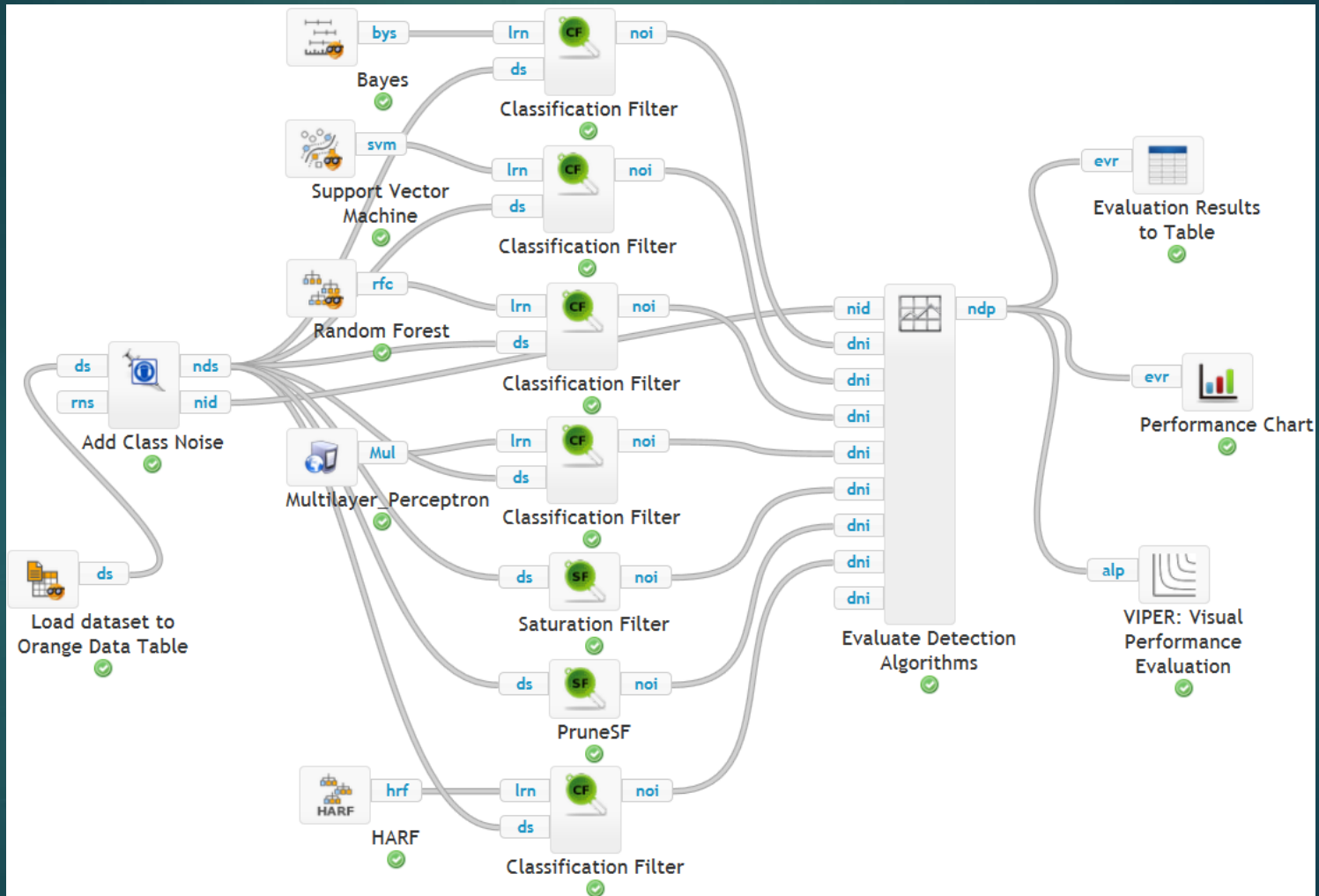
    - Is not random

# Detecting noise and outliers

- ▶ Errors and exceptions are:

  - ▶ Inconsistencies with common patterns

  - ▶ Great deviations from expected values

  - ▶ Hard to describe

# Detecting noise and outliers

▶ Identify the "laws" of the data

▶ Build models
  ▶ Patterns and rules = "laws" of the data

▶ **Errors and exceptions**
  ▶ **Do NOT obey the laws (models)**

# Classification noise filtering

► Model the data

► What can't be modeled is considered noise



► Can use any learning algorithm

(Brodley & Friedl 1999)

# Example Workflow

# Ensembles



- Combine predictions of various models

- To overcome weaknesses or bias of individual models

- Averaging, Majority voting, Consensus voting, Ranking, etc.

# Example Workflows

▶ Ensembles of noise filters

# Example Workflows

- NoiseRank

# Try it out

- Noise filtering using ensembles (with performance evaluation)
  - http://clowdflows.org/workflow/245/

- NoiseRank
  - http://clowdflows.org/workflow/115/

- Clowdflows:
  - Noise Handling
  - Orange, Weka classification
  - Performance evaluation

- Need help or advice: borut.sluban@ijs.si