

Chapter Summary: Problem Definition and Foundations of Data Mining

NOVA IMS Course Notes

Introduction

This chapter introduces the fundamental aspects of problem definition in Data Mining, contrasting traditional statistical approaches with modern data science paradigms. It emphasizes the importance of data-driven methods and the challenges posed by large, complex datasets.

1 Overview

Classification is a fundamental task in supervised machine learning. It involves assigning labels or categories to data instances based on their characteristics. Classifiers are the models that perform this task, using features derived from the data to make predictions.

2 Key Terms: Classifiers, Classes, Features, and Attributes

Term	Meaning
Classifier	A model that decides which category an example belongs to.
Classification	The task of assigning class labels to examples.
Classes	The possible groups or labels the examples belong to.
Features	The information used by the model to make decisions.
Attributes	The data columns or properties describing each example (often used as or turned into features).

Note: In databases and data warehouses, we often say *attributes*, while in machine learning we say *features*. They generally refer to the same underlying information but at different stages: attributes when stored, features when used for modeling.

3 Classifiers and Supervised Learning

A classifier is typically used in **supervised learning**, where the training data comes with known labels. The model learns to map features to classes. For example, a fruit classifier might use the color and weight of a fruit to predict whether it is an apple or an orange.

3.1 Formal Definition

In data science terms, a classifier is the trained model that performs classification — i.e., assigns a label (class) to an input example based on its features. Formally, a classifier can be represented as:

$$\text{Classifier: } f(\mathbf{x}) \rightarrow y$$

where: \mathbf{x} = input features (e.g., weight, color), y = predicted class (e.g., “apple” or “orange”).

3.2 Other Learning Paradigms

- **Unsupervised learning:** No labels are available; models find patterns or clusters in the data. Classifiers in the traditional sense are not used here.
- **Semi-supervised learning:** Some labeled data is available, but most of the data is unlabeled. Classifiers can leverage both labeled and unlabeled data.
- **Self-supervised learning:** Labels are generated from the data itself (pseudo-labels), which allows classifiers to be trained without human-provided labels.

Summary: Classical classifiers are mainly for supervised learning because they require labels. However, modern techniques extend classification concepts to scenarios where labels are partially available or generated automatically.

Key Themes

1. Data-Centric Approach

- Emphasizes the power of data over complex theoretical models.
- Advocates for using **simple models with large datasets** rather than complex models with limited data.
- Encourages the use of **unsupervised learning** and **non-parametric models** to let data “speak for itself.”

2. Statistics vs. Data Science

- Traditional statistics: small, clean, static, i.i.d. data.
- Data mining: large, messy, dynamic, non-i.i.d. data.
- Shift from **statistical significance** to **substantive significance**.

3. Challenges in Data Mining

- **Curse of Dimensionality**: Difficulty in high-dimensional spaces.
- **Nonstationarity**: Changing data distributions over time.
- **Selection Bias**: Incomplete data due to past decisions.
- **Spurious Correlations**: False patterns from large data searches.

4. Causality vs. Correlation

- Highlights the danger of confounding variables.
- Stresses the need for **causal relationships** between inputs and outputs.
- Uses examples (e.g., ice cream sales vs. drowning) to illustrate spurious correlations.

5. Data Types and Input Space

- Variable types: nominal, ordinal, interval, ratio.
- Importance of covering the **input space** for model reliability.
- Challenges of **extrapolation** vs. **interpolation**.
 - • Interpolation involves predicting a value inside the domain and/or range of the data.
 - • Extrapolation involves predicting a value outside the domain and/or range of the data

6. Model Evaluation

- Introduces concepts of **separability** and **Bayes error**.
- Discusses the importance of data coverage and avoiding unknown areas in prediction.

Conclusion

The chapter sets the foundation for data mining by highlighting the importance of a well-defined problem, a data-driven mindset, and an awareness of the pitfalls associated with large-scale data analysis. It bridges traditional statistical thinking with modern computational approaches, preparing the reader for the practical challenges of real-world data mining.