

STA 3180 Statistical Modelling: Data Mining

STA 3180 Statistical Modelling - Data Mining Lecture Notes

Introduction

Data mining is the process of discovering patterns in large datasets. It involves the use of algorithms and statistical models to uncover hidden relationships between variables and to identify trends and anomalies in data. Data mining can be used to make predictions about future events, to identify customer segments, and to develop marketing strategies.

Key Concepts

- **Data Exploration**: Data exploration is the process of examining a dataset to identify patterns, trends, and relationships. It involves visualizing data, performing statistical tests, and using machine learning algorithms to uncover insights.
- **Data Cleaning**: Data cleaning is the process of removing or correcting errors and inconsistencies in a dataset. It involves identifying and correcting errors, filling in missing values, and transforming variables.
- **Data Transformation**: Data transformation is the process of converting data from one form to another. It involves converting categorical variables into numerical variables, scaling variables, and normalizing data.
- **Data Modeling**: Data modeling is the process of creating a mathematical model to represent a dataset. It involves selecting an appropriate model, fitting the model to the data, and evaluating the model's performance.
- **Data Visualization**: Data visualization is the process of creating visual representations of data. It involves creating charts, graphs, and maps to make data easier to understand.

Definitions

- **Algorithm**: An algorithm is a set of instructions that can be used to solve a problem. Algorithms are used in data mining to identify patterns and relationships in data.
- **Statistical Model**: A statistical model is a mathematical representation of a dataset. It is used to make predictions about future events and to identify trends and anomalies in data.
- **Data Mining**: Data mining is the process of discovering patterns in large datasets. It involves the use of algorithms and statistical models to uncover hidden relationships between variables and to identify trends and anomalies in data.

- **Machine Learning**: Machine learning is the process of using algorithms to identify patterns and relationships in data. It is used to make predictions about future events and to identify customer segments.

Practice Multiple Choice Questions

Q1. What is data mining?

- A. The process of discovering patterns in large datasets.
- B. The process of creating visual representations of data.
- C. The process of converting data from one form to another.
- D. The process of using algorithms to identify patterns and relationships in data.

Answer: A. The process of discovering patterns in large datasets.