

# Data Mining - Lecture Notes

Maastricht University

October 16, 2025

## Contents

<b>1</b>	<b>Getting to Know the Data</b>	<b>1</b>
<b>2</b>	<b>Data Preprocessing</b>	<b>2</b>
<b>3</b>	<b>Regression</b>	<b>4</b>
<b>4</b>	<b>Classification</b>	<b>7</b>

# 1 Getting to Know the Data

## Lecture 2: Getting to Know the Data

### Data

Data objects are described by variables. A variable  $V$  represents a property or characteristic of an object that may vary, either from one object to another or from one time to another.

### Variable Definition

A variable  $V$  is a quadruple  $\langle \text{Name}, \text{Domain}, \text{Operations}, \text{Scale} \rangle$ :

- **Name:** The name of  $V$
- **Domain:** The set of values of  $V$
- **Operations:** The set of operations allowed over Domain
- **Scale:** A rule that associates a value from Domain for the variable  $V$  when it represents an object  $o$ .

### Types of Data

- Record Based Data (Transactions)
- Graph Based Data (WWW)
- Ordered Data (Genomics)

### Measuring Data

#### Mean

- Population Mean:  $\mu = \frac{1}{N} \sum_{i=1}^N x_i$
- Sample Mean:  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
- Weighted Mean:  $\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$

#### Median

Middle value if odd number of values, average of middle two values if even number of values

#### Mode

Value that occurs most frequently

### Symmetric vs Skewed Data

- **Symmetric:** Data is symmetrically distributed around the mean. The mean, median, and mode are all equal.
- **Skewed:** Data is not symmetrically distributed around the mean.
  - **Left Skewed (Negative):** Mean < Median < Mode
  - **Right Skewed (Positive):** Mode < Median < Mean

## Dispersion

- **Range:**  $\max(x) - \min(x)$
- **Quantile:** At most  $n(k/q)$  values will be smaller
- **IQR:**  $Q_3 - Q_1$  (middle 50% of data)
- **Outliers:** Values outside  $[Q_1 - 1.5 \times \text{IQR}, Q_3 + 1.5 \times \text{IQR}]$
- **Variance:**  $\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$
- **Standard Deviation:**  $\sigma = \sqrt{\sigma^2}$

## Data Characteristics

- **Dimensionality:** Number of variables/features. The curse of dimensionality refers to the exponential increase in data required to densely populate space as the dimension increases.
- **Sparsity:** Proportion of missing/zero values in the data.
- **Resolution:** Level of detail or aggregation in the data.

## Normal Distribution

- Bell-shaped curve
- $\mu = \bar{x}$
- $\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$
- **68-95-99.7 Rule:**
  - 68% within  $\mu \pm \sigma$
  - 95% within  $\mu \pm 2\sigma$
  - 99.7% within  $\mu \pm 3\sigma$

## Variable Types

- **Nominal:** Categories without order (gender, color, zip code)
- **Ordinal:** Ordered categories (education level, income level)
- **Interval:** Ordered with equal intervals (temperature in °C, dates)
- **Ratio:** Interval with true zero (height, weight, age)

## Statistical Plots

- **Boxplot:** Five-number summary (min, Q1, median, Q3, max) and outliers
- **Histogram:** Shows frequency distribution of numerical data
- **Quantile Plot:** Plots data against theoretical quantiles (index  $f = \frac{i-0.5}{n}$ )
- **Q-Q Plot:** Compares two distributions using their quantiles
- **Scatter Plot:** Shows relationship between two numerical variables

## Outlier Handling

1. Remove if erroneous
2. Transform (log, square root)
3. Use robust statistics (median, IQR)
4. Cap/floor extreme values

## 2 Data Preprocessing

## Lecture 3: Data Preprocessing

### Overview

The detection and correction of data quality problems. The use of algorithms that can tolerate poor data quality. Data can be:

- Inconsistent: data transformations, technology problems, human errors
- Incomplete: missing values, incomplete records
- Inaccurate: errors in data entry, data transformations, technology problems
- Outdated: data transformations, technology problems

### Data Cleaning

Converting data so that it becomes consistent, complete, accurate, and up-to-date. It's realized by filling missing values, removing duplicates, smoothing noise, and resolving inconsistencies.

### Handling Noisy Data

- Clustering (detect and remove outliers)
- Computer and human inspection

### Handling Missing Data

- Filling manually
- Using the variable mode, median, or mean

### Data Integration

Combining data from different sources.

### Possible Problems

- Different variables have the same name
- Similar variables have different names
- Redundant variable: can be detected with Chi-Square, Covariance analysis

### Chi-square Test

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

### Covariance

$$\text{Cov}(X, Y) = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)$$

### Data Reduction

Obtaining a reduced set of variables that are sufficient to represent the data.

## Strategies

- **Principal Component Analysis (PCA):** Converts variables into a new set of variables that are uncorrelated and capture the maximum variance.
- **Multidimensional Scaling (MDS):** Finds a low-dimensional representation that preserves pairwise distances (or dissimilarities) between points.
- **Feature Selection:** Select a subset of variables that are most relevant to the task.
- **Clustering:** Group similar objects together.
- **Sampling:** Main strategy for data reduction in data mining. The sample must be representative of the population.
  - Without replacement: each object is selected only once
  - With replacement: each object can be selected multiple times
  - Stratified sampling: data is split into partitions and a sample is taken from each partition

## Data Valuation

Seeks to assign a numerical value to an individual's data in the trade of data. The issue is the time and cost of data valuation. Complexity is above  $O(2^N)$ .

## Data Transformation and Discretization

A function that maps the entire set of values of a given variable to a new set of replacement values.

### Methods

- **Normalization:** Scales the values to a range, such as  $[0, 1]$  or  $[-1, 1]$ .
- **Smoothing:** Reduces noise in the data.
- **Variable/Feature Construction:** Creates new variables from existing variables.

### Normalization

- **Min-Max:** Scales the values to a range, such as  $[0, 1]$  or  $[\text{New Min}, \text{New Max}]$

$$\bar{x} = \frac{x - \min(x)}{\max(x) - \min(x)}$$

- **Z-Score:** Scales the values to have a mean of 0 and a standard deviation of 1.

$$\bar{x} = \frac{x - \mu}{\sigma}$$

### Discretization

Divides the range of continuous values into a set of intervals. The intervals are called bins and can replace the original values. Clustering can also be used to find the intervals.

- **Binning:** Can be done with equal width or equal frequency (depth)

## 3 Regression

## Lecture 4: Regression

### Regression

Given a set of variables  $X$ , we want to predict a target variable  $Y$ . There's an unknown function  $f$  that maps  $X$  to  $Y$ .

We assume that  $Y = f(X) + \epsilon$  where  $\epsilon$  is the noise.  $f(X)$  is an optimal function that minimizes the error. The random error term  $\epsilon$  is assumed to be independent of  $X$  and has a mean of 0 and cannot be reduced by any model.

The optimal function  $f(X)$  is the one that minimizes the error. The error is defined as the difference between the predicted value and the actual value.

### Parametric Methods

Parametric methods assume a specific form for the function  $f(X)$ . They involve two steps:

1. We assume a form for  $f(X)$ .

$$f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

2. Using training data, we find the parameters of  $f(X)$  that minimize the error.

Because of the restricted flexibility of parametric methods, they are open to **underfitting**.

- **Underfitting:** The model is too simple to capture the relationship between  $X$  and  $Y$ .

### Non-Parametric Methods

Non-parametric methods do not assume a specific form for the function  $f(X)$ . They are flexible and can capture complex relationships between  $X$  and  $Y$ , but they require more data to estimate the function. They are open to **overfitting**.

- **Overfitting:** The model corresponds too closely to the training data and fails to generalize to new data.

### Assessing Model Performance

- **Mean Squared Error (MSE):**

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

### Bias-Variance Tradeoff

The expected error of a model is the sum of three components:

- **Bias:** The error due to the model's inability to capture the relationship between  $X$  and  $Y$ .
- **Variance:** The error due to the model's sensitivity to the training data.
- **Irreducible Error:** The error due to the noise in the data.

$$\text{Expected Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

### Linear Regression

The linear regression model assumes a linear relationship between  $X$  and  $Y$ . We assume a true dependency:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$$

## Assessing the Model

- **Residuals:**  $y_i - \hat{y}_i$
- **Residual Sum of Squares (RSS):**  $\sum_{i=1}^N (y_i - \hat{y}_i)^2$
- **Residual Standard Error (RSE):**  $\sqrt{\frac{\text{RSS}}{N-2}}$
- **R-squared:**  $R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$
- **Total Sum of Squares (TSS):**  $\text{TSS} = \sum_{i=1}^N (y_i - \bar{y})^2$

RSE is the standard deviation of the residuals, or the irreducible error epsilon.

## Additive Assumption

The effect of any variable is independent of the values of the effects of other variables.

To avoid the additive assumption, we can use interaction terms. We move from the model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

To the model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_{12} X_1 X_2 + \epsilon$$

The interaction term introduces a non-linear relationship between  $X_1$  and  $X_2$ .

## Shrinkage Methods

Shrinkage methods reduce the variance of the model by shrinking the coefficients of the variables.

### Ridge Regression

Ridge regression adds a penalty term to the RSS:

$$\text{RSS} + \lambda \sum_{i=1}^p \beta_i^2$$

Where  $\lambda$  is the regularization parameter.

The Ridge Regression model is:

$$\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$$

### Lasso Regression

Lasso regression adds a penalty term to the RSS:

$$\text{RSS} + \lambda \sum_{i=1}^p |\beta_i|$$

This has the effect of forcing some of the coefficients to zero. So the Lasso performs feature selection.

## KNN Regression

KNN regression is a non-parametric method that uses the k-nearest neighbors to predict the value of a new point.

The KNN model is:

$$\hat{y} = \frac{1}{k} \sum_{i=1}^k y_i$$

The value of k is a hyperparameter that needs to be tuned. If k is small, then the regression model is more flexible.

However, this regression is not so good in high dimensions.

## Decision Tree Regression

The decision tree regression model is a non-parametric method that uses a tree to predict the value of a new point.

Regression tree is learned by minimizing the RSS:

$$\text{RSS} = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_j)^2$$

Where  $R_j$  is the number of regions in the tree.

Regression Trees are sensitive to overfitting. To avoid this, we can use pruning. There are two approaches:

1. **Pre-pruning:** Stop the tree before it overfits.
2. **Post-pruning:** Prune the tree after it overfits.

Model Trees provide different values for different instances due to the regression models in the leaves.

## 4 Classification

### Lecture 5: Classification

#### Classification

##### Naïve Bayes

This is a parametric method that uses Bayes' theorem to predict the class of a new point.

It is based on the Bayes' theorem:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

It assumes that the features are independent given the class and are identically distributed. Independent samples are samples that are drawn from the same distribution.

There are two types of classifiers:

- **Discrete classifiers:** assign a class label to a test instance
- **Score classifiers:** assign a continuous score for each class and can be assigned to a test instance.

An optimal Bayes rule assumes knowledge of:

- The prior distribution  $p(y)$
- The distribution  $p(x|y)$  for each class  $y$
- The distribution  $p(x)$

Naïve Bayes classifier naively assumes that the input variables are conditionally independent. The function is:

Naïve Bayes can be a linear or non-linear method for classification depending on the properties of the input variables.

Naïve Bayes does not have any parameter to control the bias-variance trade-off. The only way can be explicit feature selection.



## Logistic Regression

Instead of predicting  $Y$ , we predict  $P(Y = 1|X)$  using the logistic (sigmoid) function:

$$P(Y = 1|X) = \frac{e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p}}$$

For more than two classes we use the softmax function:

$$P(Y = k|X) = \frac{e^{\beta_{k0} + \beta_{k1} X_1 + \beta_{k2} X_2 + \dots + \beta_{kp} X_p}}{\sum_{i=1}^K e^{\beta_{i0} + \beta_{i1} X_1 + \beta_{i2} X_2 + \dots + \beta_{ip} X_p}}$$

We can estimate the parameter  $\beta$  using maximum the likelihood function:

$$l(\beta_0, \beta) = \prod_{i=1}^n P(Y = y_i | x_i)$$

Where if  $Y = 1$ :

$$P(Y = y_i | x_i) = \frac{e^{\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}}}{1 + e^{\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}}}$$

And if  $Y = 0$ :

$$P(Y = y_i | x_i) = \frac{1}{1 + e^{\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}}}$$

### Summary

- Logistic regression is a parametric method for classification.
- Logistic regression is a linear method for classification.
- Logistic regression estimates class probabilities. It does not make a classification decision; i.e., it is a scoring classifier.
- The variance of logistic regression can be reduced using shrinkage methods based on ridge regression (ridge logistic regression).

## Support Vector Machine

Support Vector Machines (SVMs) approach the two-class classification problem in a direct way. SVM tries indeed to separate the classes in instance space  $X$ .

Among all separating hyperplanes, find the one that makes the biggest gap or margin between the two classes.

To convert SVM to a scoring classifier, we can use the decision function:

$$f(x) = \frac{1}{1 + e^{Af(x)+B}}$$

where  $A$  and  $B$  are parameters that can be estimated using maximum likelihood.

Higher  $C$  values imply low flexibility (high bias, low variance). Lower  $C$  values imply high flexibility (low bias, high variance).

### Summary

- SVM is a discrete classifier. It provides a classification (no probability)!
- SVM can be converted to a scoring classifier using signed distance to hyperplane (directly or using the Platt scaling).
- SVM is a parametric method for binary classification.
- SVM is a linear method for classification.
- SVMs handle nonseparability problems using: Soft-margins and Kernels.

**Feature Expansion** Enlarge the feature space  $X$  by adding new features:  $X^2$ ,  $X^3$ ,  $X_1 X_2$ , ... Fit a support-vector classifier in the enlarged space. This results in non-linear decision boundaries in the original space.

**Kernel Support Vector Machine** We have a kernel function  $K(x_i, x_j)$  that computes the similarity between  $x_i$  and  $x_j$ . Then the decision function is:

$$f(x) = \beta_0 + \sum_{i=1}^n \alpha_i K(x_i, x)$$

## Decision Trees

Each interior node tests a variable. Each branch corresponds to a variable value. Each leaf node is labeled with a class (class node).

```
function Classify(x: instance, node: variable containing a node of DT)
  if node is a classification node then
    return the class of node
  else
    determine the child of node that matches x
    return Classify(x, child)
  end if
end function
```

It is okay for the training data to contain missing values. Decision trees can be used even if instances have missing variables.

## Basic Algorithm:

1.  $X \leftarrow$  the "best" decision variable for a node  $N$ .
2. Assign  $X$  as decision variable for the node  $N$ .
3. For each value of  $X$ , create new descendant of the node  $N$ .
4. Sort training instances to leaf nodes.
5. IF training examples perfectly classified, THEN STOP. ELSE iterate over new leaf nodes.

**Entropy** Let  $S$  be a sample of training examples, and  $p_+$  is the proportion of positive examples in  $S$  and  $p_-$  is the proportion of negative examples in  $S$ . Then: entropy measures the impurity of  $S$

$$H(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

**Bias-Variance Tradeoff** Decision trees have in general high variance.

- The bias of decision trees decreases with the size of the trees.
- The variance of decision trees increases with the size of the trees.

## Overfitting

- **Pre-pruning:** stop growing the tree earlier, before it reaches the point where it perfectly classifies the training data.
- **Post-pruning:** Allow the tree to overfit the data, and then post-prune the tree.

Validation set is a set of instances used to evaluate the utility of nodes in decision trees. The validation set has to be chosen so that it is unlikely to suffer from same errors or fluctuations as the training set.

## Summary

- DTs are discrete classifiers. They can estimate probabilities by normalizing class scores in each leaf node.
- Decision Trees (DT) for a non-parametric method for classification.
- DTs are a non-linear method for classification.

## Decision Rules

Decision rules are rules with the following form:

if {conditions} then concept C

### Summary

- Decision Rules (DRs) are discrete classifiers. They can estimate probabilities by normalizing class scores in each rule.
- Decision Rules form a non-parametric method for classification.
- DRs are a non-linear method for classification.
- DRs are usually simpler than decision trees on the same data.

## K-NN Classification

k-NN Classifier is a non-parametric classifier. To estimate a class value  $y$  for a given test instance  $x$ , find a set  $NN$  of the  $k$  closest instances to  $x$  in training data  $Tr$ .

- **Discrete Classification:** output the majority class among the instances in  $NN$ .
- **Scoring Classification:** output the score for each class among the instances in  $NN$ . If the scores are normalized we estimate class probabilities.

The value of  $k$  controls the flexibility of the k-NN classifier. The smaller that value the more flexible is the k-NN classifier (the higher the variance and lower the bias).

### Notes:

- Continuous variables should be normalized. Otherwise, the variables with bigger domains prevail!
- Discrete variables do not pose problems since distances are based on value matches.

### Advantages

1. The NN classifier can estimate complex class borders locally and differently for each new test instance.
2. The NN classifier provides good generalization performance on many domains.
3. The NN classifier learns very quickly.
4. The NN classifier is robust to noisy training data.
5. The NN classifier is intuitive and easy to understand which facilitates implementation and modification.

### Disadvantages

1. The NN classifier has large storage requirements because it has to store all the data.
2. The NN classifier is slow during instance classification because all the training instances have to be visited.
3. The generalization performance of the NN classifier degrades with increase of noise in the training data.
4. The generalization performance of the NN classifier degrades with increase of irrelevant variables.

### Summary

- Nearest-Neighbor (NN) Classifier is a non-parametric method for classification.
- NN is a non-linear method for classification.
- NN can be a discrete classifier and a scoring classifier depending on how we handle the class statistics of the nearest neighbors.
- The bias-variance trade-off can be controlled by the parameter  $k$ .