

Chapter 3: Classification Techniques

Your Name

Your Institution

July 19, 2025

Overview of Classification Techniques

Classification is a fundamental technique in data mining involving the assignment of items to target categories or classes. It is essential for organizing and analyzing large amounts of data, enabling easier decision-making.

Importance in Data Categorization

- **Data Organization:** Systematically categorizes data for easier retrieval and analysis.
- **Predictive Insights:** Enables predictions about future behavior and outcomes, such as spam detection in emails.
- **Decision Support:** Assists in decision-making across various fields like finance, healthcare, and marketing.

- 1 **Supervised Learning:** Classification typically occurs in supervised settings using labeled data.
- 2 **Features and Classes:** Each data point is represented by features (attributes) used to predict the class (category).

Example

In email classification, features could include word frequencies, while classes might be 'spam' or 'not spam'.

Objectives of the Chapter

- ❶ **Understanding Classification:** Define classification and its role in data mining.
- ❷ **Techniques Overview:** Explore techniques like:
 - Decision Trees
 - Support Vector Machines (SVM)
 - Naive Bayes
 - Neural Networks
- ❸ **Evaluation Metrics:** Learn about metrics like accuracy, precision, recall, and F1-score.
- ❹ **Real-world Applications:** Discuss applications in various industries, showcasing their value.

Formula and Concepts to Note

Confusion Matrix

	Predicted Positive	Predicted Negative
Actual Positive	TP	FN
Actual Negative	FP	TN

Where:

- TP = True Positives
- TN = True Negatives
- FP = False Positives
- FN = False Negatives

Performance Metrics

- Accuracy = $\frac{TP+TN}{TP+TN+FP+FN}$
- Precision = $\frac{TP}{TP+FP}$
- Recall (Sensitivity) = $\frac{TP}{TP+FN}$

Summary

Classification techniques are vital for transforming raw data into actionable insights. Understanding these concepts sets the foundation for effective data analysis and predictive modeling, which will be explored in greater detail throughout this chapter. Let's delve deeper into what classification entails in the next slide!

What is Classification?

Definition of Classification

Classification is a supervised machine learning technique used in data mining to systematically organize data into predefined categories or classes. The primary goal is to predict the categorical label of new, unseen instances based on past observations and knowledge about the data.

Key Concepts of Classification

- **Supervised Learning:** Refers to an algorithm learning from labeled training data.
- **Example:** In email filtering, 'spam' and 'not spam' are the predefined classes.

Example of Classification

Consider classifying emails as "spam" or "not spam":

- 1 Collect a dataset of emails with labels (Spam or Not Spam).
- 2 Use a classification algorithm (e.g., Naive Bayes) to learn the characteristics that differentiate spam from non-spam.
- 3 Test the model with new emails; it will classify them based on learned characteristics.

Classification vs. Clustering

- **Labeling:**

- Classification: Involves labeling data with predefined categories.
- Clustering: Groups data without predefined labels to identify patterns.

- **Purpose:**

- Classification: Predicts categories of new instances based on learned patterns.
- Clustering: Discovers inherent groupings in data without supervision.

Illustration of Difference

- **Classification Example:** Predicting whether a patient has a disease (yes/no) using medical records.
- **Clustering Example:** Grouping customers based on purchasing behavior without defined categories.

Key Points to Emphasize

- Classification requires labeled data for model training.
- Focuses on prediction as opposed to clustering's discovery of inherent groupings.
- Effectiveness relies on the quality of training data and relevance of features.

Conclusion

Classification plays a pivotal role in predictive analytics, enabling businesses and researchers to make informed decisions based on historical data and observed trends.

Types of Classification Techniques

Overview of Classification Techniques

Classification is a crucial task in data mining where we categorize data into predefined classes. Here, we will examine four major classification techniques:

- Decision Trees
- Support Vector Machines (SVM)
- K-Nearest Neighbors (KNN)
- Neural Networks

1. Decision Trees

- **Description:** A flowchart-like structure representing decisions based on attributes.
- **Example:** Classifying loan suitability based on income, credit score, and employment status.
- **Key Point:** Intuitive and easy to interpret; prone to overfitting unless properly pruned.

2. Support Vector Machines (SVM)

- **Description:** A supervised learning model that identifies the hyperplane that separates classes in high-dimensional space.
- **Example:** Finding the best line to separate 2D points of cats and dogs.
- **Key Point:** Effective in high-dimensional spaces and robust against overfitting, especially with the right kernel function.

3. K-Nearest Neighbors (KNN)

- **Description:** An instance-based learning algorithm that classifies a point based on its neighbors.
- **Example:** With $K=3$, if two neighbors are "A" and one is "B," the new point is classified as "A."
- **Key Point:** Easy to understand and implement, but computationally expensive for large datasets and sensitive to irrelevant features.

4. Neural Networks

- **Description:** Consists of interconnected layers of nodes (neurons) that learn by adjusting weights based on input data.
- **Example:** Classifying images, such as identifying handwritten digits.
- **Key Point:** Highly flexible and can model complex relationships, but requires more data and computational power.

Summary and Conclusion

Summary

Each classification technique comes with its strengths and weaknesses. The choice depends on:

- The nature of the data
- The specific problem to be solved
- Project resource constraints

Understanding these techniques helps in selecting the most appropriate method.

Conclusion

Selecting the right classification technique is crucial for building accurate and reliable predictive models.

Definition

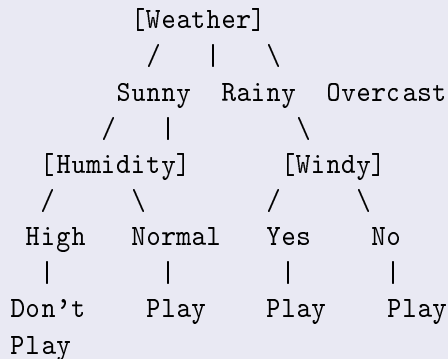
A decision tree is a flowchart-like structure used for decision-making and predictive modeling. It helps classify data by creating a model that predicts the value of a target variable based on several input features.

Decision Trees - Structure

- **Components:**

- **Node:** Represents a feature or attribute (e.g., weather, age).
- **Branch:** Represents the decision rule (e.g., 'Is the weather sunny?').
- **Leaf Node:** Represents the outcome or class label (e.g., 'Play' or 'Don't Play').

Example Structure



Decision Trees - How They Work

- **Splitting:** At each node, the dataset is split into subsets based on a feature to create homogeneous subsets.
- **Criteria for Splitting:**
 - **Gini Impurity:** Measures the likelihood of a random instance being incorrectly labeled.
 - **Entropy:** Measures the randomness in the dataset. The goal is to reduce entropy with each split.

Formulas

$$Gini = 1 - \sum (p_i^2) \quad (1)$$

$$Entropy = - \sum p_i \log_2(p_i) \quad (2)$$

Decision Trees - Advantages and Disadvantages

- **Advantages:**

- Intuitive and easy to interpret.
- No need for data normalization.
- Handles both numerical and categorical data.

- **Disadvantages:**

- Overfitting: Can create overly complex trees.
- Instability: Small changes in data can lead to different trees.

- **Industry Use Cases:**

- **Finance:** Credit scoring.
- **Healthcare:** Diagnosis classification.
- **Marketing:** Customer segmentation.

Decision Trees - Key Takeaways

- Decision trees are versatile tools for classification tasks.
- They offer an intuitive method for analysis and predictive modeling.
- Understanding their structure and principles is crucial for real-world applications.

Decision Trees - Conclusion

Decision trees are an essential part of classification techniques, bridging the gap between data analysis and decision-making processes. In the next slide, we will explore Support Vector Machines (SVM) and their application in classification tasks.

Support Vector Machines (SVM) - Overview

Key Concept

Support Vector Machines (SVM) are powerful supervised learning models used primarily for classification and regression tasks.

The objective of SVM is to find the best hyperplane that separates data points from different classes in a high-dimensional space.

Support Vector Machines (SVM) - Working Principle

- 1 **Hyperplane Definition:** A hyperplane is a flat affine subspace dividing space into two half-spaces.
- 2 **Margin Calculation:** SVM aims to maximize the margin, defined as the distance between the hyperplane and the nearest data points (Support Vectors).
- 3 **Mathematical Representation:**

$$w \cdot x + b = 0 \quad (3)$$

Where:

- w : Weight vector
- x : Input feature vector
- b : Bias term

4 Optimization Problem:

$$\text{Minimize } \frac{1}{2} ||w||^2 \quad \text{subject to } y_i(w \cdot x_i + b) \geq 1 \quad (4)$$

Where y_i is the class label (+1 or -1).

- 5 **Kernel Trick:** Allows SVMs to handle non-linear separation by projecting data into higher dimensions. Common kernels include:
- Linear kernel
 - Polynomial kernel
 - Radial Basis Function (RBF) kernel

Support Vector Machines (SVM) - Applications and Key Points

Example Application

- Image Classification: SVM can classify images into categories (e.g. detecting cats vs. dogs).

Key Points to Emphasize

- Support Vectors are crucial for determining the hyperplane.
- SVM is effective in high-dimensional spaces with a clear margin of separation.
- Regularization parameter (C) helps prevent overfitting.

K-Nearest Neighbors (KNN) - Introduction

- KNN is a simple yet powerful classification algorithm in machine learning.
- It is instance-based, deferring computation until function evaluation.
- No model is built; predictions are made based on data point proximity.

K-Nearest Neighbors (KNN) - Algorithm Overview

- 1 Choose the number of neighbors, *K*.
- 2 Calculate distance between the new data point and existing points.
 - Euclidean Distance:
$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$$
 - Manhattan Distance: $d(p, q) = |p_1 - q_1| + |p_2 - q_2| + \dots + |p_n - q_n|$
- 3 Sort distances and identify the *K* nearest neighbors.
- 4 Assign the most common class label from the *K* neighbors.

K-Nearest Neighbors (KNN) - Factors Influencing Performance

- **Choosing the Right K:**
 - Small *K*: Sensitive to noise, risk of overfitting.
 - Large *K*: Smoother decision boundary, may overlook local patterns.
- **Distance Metric:**
 - Affects classification significantly (e.g., Euclidean vs. Hamming distance).
- **Feature Scaling:**
 - Sensitive to different scales; standardization or normalization is essential.
- **Dimensionality:**
 - "Curse of Dimensionality" can make distance measures less informative.

K-Nearest Neighbors (KNN) - Key Points

- KNN is intuitive and flexible for classification tasks.
- Performance is influenced by:
 - Choice of K
 - Distance metric
 - Feature scaling
- Data preprocessing is crucial for effective KNN application.

Overview of Neural Networks as Classifiers

Neural Networks are a powerful class of models used in machine learning and artificial intelligence, designed to recognize patterns within data. They are particularly effective for tasks involving classification where the outcome is categorical.

- **Neurons**: Basic unit, processes input data.
- **Layers**: Input layer, hidden layers, output layer.
- **Activation Functions**: Introduce non-linearity (e.g., ReLU, Sigmoid).

Neural Networks - Shallow vs. Deep Learning

Shallow Neural Networks

- **Definition**: 1 hidden layer.
- **Use Cases**: Simpler problems, linearly separable data.
- **Example**: Classifying flower types based on sepal and petal measurements.

Deep Learning Models

- **Definition**: Multiple hidden layers.
- **Use Cases**: Complex tasks (e.g., image recognition).
- **Example**: CNNs for image classification.

Differences Between Shallow and Deep Learning Models

Feature	Shallow Neural Networks	Deep Neural Networks
Architecture	1 hidden layer	Multiple hidden layers
Complexity	Simpler, easier to interpret	More complex, harder to interpret
Feature Learning	Manual extraction required	Learns features from data
Performance	Limited on complex tasks	Superior on large datasets
Training Time	Faster	Longer due to complexity

Key Points

- Mimics the human brain's learning process.
- Shift from shallow to deep models indicates increased complexity in data interpretation.
- Deep learning reduces the effort in feature engineering.

Example of a Simple Neural Network Model

```
import numpy as np
from keras.models import Sequential
from keras.layers import Dense

# Example: Simple Shallow Neural Network
model = Sequential()
model.add(Dense(10, input_dim=8, activation='relu'))
    # 10 neurons, input size of 8
model.add(Dense(2, activation='softmax')) # 2 output
    neurons for binary classification
model.compile(loss='categorical_crossentropy',
              optimizer='adam', metrics=['accuracy'])
```

Conclusion

Neural networks represent significant advancements in machine learning, evolving from shallow understanding to deep learning capabilities.

Model Evaluation Metrics

Evaluating classification models is crucial for determining effectiveness. Key metrics include:

- 1 Accuracy
- 2 Precision
- 3 Recall (Sensitivity)
- 4 F1 Score

1. Accuracy

Definition

Accuracy measures the overall correctness of the model. It is the ratio of correctly predicted instances ($TP + TN$) to the total instances ($TP + TN + FP + FN$).

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (5)$$

- Example: If a model makes 80 correct predictions out of 100 total predictions, accuracy is 80
- Key Point: Accuracy can be misleading, particularly in unbalanced datasets.

2. Precision

Definition

Precision measures the correctness of positive predictions. It is the ratio of true positives to the sum of true positives and false positives.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (6)$$

- Example: If a model identifies 10 instances as positive, but only 7 are truly positive, precision is $\frac{7}{10} = 0.7$ or 70
- Key Point: High precision indicates a low number of false positives, crucial in applications like spam detection.

3. Recall (Sensitivity)

Definition

Recall measures how well the model identifies actual positive instances. It is the ratio of true positives to the sum of true positives and false negatives.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (7)$$

- Example: If there are 10 actual positive instances, and the model correctly predicts 8, recall is $\frac{8}{10} = 0.8$ or 80
- Key Point: High recall signifies that many actual positives are correctly identified, crucial in fields like medical diagnosis.

4. F1 Score

Definition

The F1 Score is the harmonic mean of precision and recall, balancing the two. It's especially useful when class distribution is uneven.

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (8)$$

- Example: With precision at 70

$$F1 = 2 \times \frac{0.7 \times 0.8}{0.7 + 0.8} \approx 0.746$$

- Key Point: The F1 Score is advantageous when both false positives and false negatives matter significantly.

Conclusion

In summary, evaluating classification models involves using multiple metrics, each providing distinct insights into performance. Depending on the context, one may prioritize precision, recall, or F1 Score over accuracy to better reflect how well a model meets its objectives. Consider the specific needs of your application when choosing evaluation metrics.

Challenges in Classification - Overview

- Classification challenges include:
 - Overfitting
 - Underfitting
 - Imbalanced datasets
- Understanding these challenges is crucial for developing effective classification models.

Challenge 1: Overfitting

Definition

Overfitting occurs when a model learns the noise in the training data instead of the underlying pattern, resulting in high accuracy on training data but poor generalization to unseen data.

- **Example:** A model memorizing unique features in images (e.g., backgrounds) instead of general characteristics leads to misclassification of new images.

Key Points to Address Overfitting

- Simplify the model
- Use regularization techniques (L1, L2)
- Implement cross-validation

Challenge 2: Underfitting

Definition

Underfitting occurs when a model is too simple to capture the underlying trend of the data, leading to poor performance on both training and testing datasets.

- **Example:** A linear model fitting a highly nonlinear dataset fails to capture complexities, resulting in low accuracy.

Key Points to Address Underfitting

- Increase model complexity
- Enhance feature engineering
- Perform hyperparameter tuning

Challenge 3: Imbalanced Datasets

Definition

An imbalanced dataset has class representation that is not approximately equal, with one class dominating others (e.g., 950 negative cases vs. 50 positive cases).

- **Consequences:** Classifiers may become biased towards the majority class, affecting performance on minority classes.

Strategies to Handle Imbalanced Datasets

- Resampling Techniques
 - Oversampling (e.g., SMOTE)
 - Undersampling
- Cost-sensitive Learning
- Specialized Algorithms (e.g., decision trees, ensemble methods)

Summary and Techniques

- Classification challenges significantly affect model performance.
- Addressing overfitting and underfitting, along with managing imbalanced datasets, is crucial for robust model development.

Regularization (L2 Penalty Example)

Loss function with L2 regularization:

$$\text{Loss} = \text{Loss}_{\text{original}} + \lambda \sum_{i=1}^n \theta_i^2 \quad (9)$$

where λ is the regularization parameter.

Oversampling with SMOTE (Python Code)

```
from imblearn.over_sampling import SMOTE
sm = SMOTE(random_state=42)
X_res, y_res = sm.fit_resample(X, y)
```

Conclusion and Future Trends - Overview of Classification Techniques

- We've explored various classification techniques fundamental to machine learning.
- Key methods include:
 - 1 Linear Classifiers (e.g., Logistic Regression)
 - 2 Decision Trees
 - 3 Support Vector Machines (SVM)
 - 4 Ensemble Methods (e.g., Random Forest)

• Linear Classifiers

- Use a linear combination of features.
- Example: Email filtering (spam vs. legitimate).

• Decision Trees

- Tree-like model for classification based on feature values.
- Example: Classifying patient risks.

• Support Vector Machines (SVM)

- Find optimal hyperplane for class separation.
- Example: Classifying images of pets.

• Ensemble Methods

- Combine multiple classifiers for improved accuracy.
- Example: Credit scoring models.

Conclusion and Future Trends - Future Directions

- **Deep Learning Approaches**

- CNNs and RNNs for image and speech recognition.
- Example: Autonomous vehicles.

- **Transfer Learning**

- Use pre-trained models for related tasks.
- Example: Medical image classification.

- **Explainable AI (XAI)**

- Increasing demand for model transparency.
- Example: SHAP values for model outputs.

- **Automated Machine Learning (AutoML)**

- Streamlines model selection and tuning.
- Example: Google Cloud AutoML for classification.

- **Handling Imbalanced Datasets**

- Techniques like SMOTE for managing class distribution.
- Example: Rare condition classification.

Conclusion and Future Trends - Conclusion

Classification techniques are essential for various industries including finance, healthcare, and technology. As advancements in computational capabilities continue, we anticipate novel approaches will significantly enhance the effectiveness of classification tasks in diverse applications.