

Knowledge Discovery in Databases
SS 2015

Übungsblatt 1: Aufgabenstellungen im Data Mining

Aufgabe 1-1 Data Mining Aufgaben

Welche Aufgaben für das Data Mining (Clustering, Outlier Detection, Klassifikation, etc.) verbergen sich hinter den folgenden Anwendungen? Ist die Aufgabe überwacht (supervised) oder nicht überwacht (unsupervised)?

(a) **Texterkennung/OCR:**

Beim Passieren der Brennerautobahn existiert seit einigen Jahren die Möglichkeit per E-Maut zu zahlen. Dabei wird bei Zahlungseingang das Nummernschild des Autos registriert. Beim Passieren der Mautstation fährt das Auto dann durch eine gesonderte Schranke die nur aufgeht, wenn das Nummernschild des Fahrzeugs als registriert erkannt wurde. Die Erkennung erfolgt dabei voll automatisch per digitaler Kamera.

(b) **Computer Aided Diagnosis:**

Patienten, die an Blutkrebs leiden, können in zwei Kategorien (ALL und AML) eingeteilt werden. Da sich die Therapien dieser beiden Arten teilweise sehr stark unterscheiden und sogar manchmal die Therapie für AML sehr schädlich für ALL-Patienten sein kann (und umgekehrt), versucht man, neue Patienten anhand von speziellen Daten (sog. Gen-Expressionsdaten) zu unterscheiden. Dazu werden die Daten der neuen Patienten mit den Daten der Patienten, deren Blutkrebstyp bereits bekannt ist, verglichen.

(c) **Cheat Detection**

Der Betreiber eines Multiplayer-Online-Spiels will sein System gegen verschiedene Verstöße der Benutzerrichtlinien abdecken. Dazu gehören die Verwendung von Bot-Programmen, das Manipulieren von Zeitstempeln im Kommunikation Protokoll und die Vorhersage verwendeter Zufallszahlen. Zur Erkennung von verdächtigem Verhalten wird Data Mining auf den erhältlichen Benutzerdaten verwendet.

(d) **Mensch und Maschine**

Moderne WWW-Suchmaschinen beantworten Benutzeranfragen, die aus nur einem oder wenigen Suchtermen bestehen. In der Regel liefert eine Anfrage dabei eine sehr große Ergebnismenge, die mit Hilfe eines Ranking Algorithmus nach Relevanz sortiert wird. Durch diese Sortierung kann der User dann selber entscheiden, wieviele Links er besuchen will. Die Problematik hierbei ist zum einen den Inhalt einer Ergebnisseite richtig zu erkennen. Zum anderen besteht die Notwendigkeit, dass wirklich hilfreiche Seiten höher gerankt werden als weniger hilfreiche Seiten, auch wenn beide inhaltlich zum Suchbegriff passen. Wortmehrfachheiten stellen dabei ebenfalls ein großes Problem dar. Zum Beispiel kann sich die Suche nach dem Begriff "Golf" auf das Auto, den Sport oder den geographischen Begriff beziehen. Data Mining Techniken werden hier eingesetzt um das Ranking zu optimieren und mögliche Ergebnis-Mengen nach dem jeweiligen Begriffskontext zu gruppieren.

(e) **Recommendation Systems**

Ein Online-Kaufhaus möchte für registrierte Kunden Artikel bestimmen, die dem Kunden beim Einloggen unaufgefordert angeboten werden. Dabei kann man auf die bereits gekauften Artikel des Kunden zurückgreifen, um so die Interessengebiete des Kunden besser vorhersagen zu können. Zum Beispiel

bietet es sich an, jemandem, der das Buch "Herr der Ringe" gekauft hat, auch die DVDs der Verfilmung anzubieten. Eine weitere ähnliche Aufgabe ist die Bestimmung von geeigneten Kombiangeboten zu einem bereits ausgewählten Artikel.

(f) **News Aggregation**

Eine Nachrichtenseite sammelt automatisch Meldungen aus verschiedenen Nachrichtenquellen um den Nutzer zu informieren. Da es aber häufig passiert, dass unterschiedliche Seiten über das selbe Thema berichten, sollen die Meldungen gruppiert werden. Die Überlappungen passieren dabei auf unterschiedlichen Ebenen: Zum einen gibt es natürlich breite Kategorien wie Sport und Politik, und Unterkategorien wie Fußball. Aber auch zu einem einzelnen Fußballspiel gibt es meist mehrere unterschiedliche Beiträge auf unterschiedlichen Seiten. Manche der Beiträge sind (weitgehend) identisch zu Agenturmeldungen, andere individuelle eigene Beiträge.

(g) **Extraktion von Daten / Web Scraping:**

Aus einer bekannten Filmdatenbank sollen eine Liste von Filmen und eine Liste von Schauspielern extrahiert werden (Lizenzprobleme seien für diese Aufgabe ignoriert).

(h) **Identifikation der wichtigsten Zulieferer:**

Ein großer Onlinehändler möchte wissen, welche Lieferanten für ihn am wichtigsten sind, d.h. den größten Umsatz beisteuern. Zu diesen könnten dann engere Beziehungen geknüpft werden, eine Übernahme der Firma erfolgen, oder ein neues Logistikzentrum nahe des Standortes des Lieferanten entstehen, um die Lieferzeiten zu verkürzen.

(i) **Bildsegmentation in medizinischen Bilddaten:**

Segmentation nennt man den Process des Unterteilens eines Bildes in verschiedene Teile. In der medizinischen Bildbearbeitung bedeuten diese Segmente meist verschiedene Zelltypen, Organe oder Pathologien, oder andere biologisch relevanten Strukturen. Medizinische Bildverarbeitung wird erschwert durch schlechte Bildqualität, niedrigen Kontrast und Rauschen oder andere Bildunklarheiten. Auch wenn es für Bilder schon viele Methoden gibt, werden und müssen diese meist noch für die Verwendung im medizinischen Bereich angepasst werden.

Begriffe in diesem Themenbereich sind unter anderem:

(i) Atlas-Based Segmentation:

Ein Experte bewertet einige Beispielfelder, anhand derer dann durch Extrapolation eine Aussage über das neue Bild gemacht wird. Dabei wird von den Trainingsdaten abstrahiert und daraus ein Modell entwickelt.

(ii) Shape-Based Segmentation:

Bei dieser Methode werden meist parametrisierte Modelle von Formen verwendet, die sich auf besondere Strukturen und Verläufe beziehen. Dabei wird die Form verändert um dem neuen Bild zu gleichen

(iii) Interactive Segmentation:

Ein Arzt gibt während der Operation Informationen, wie zum Beispiel die Region oder die Grenze zu einem Segment. Ein Algorithmus kann dann die Zwischenergebnisse verfeinern und somit genauer die Ausmaße eines Zelltypus definieren.