

Begriffsdefinitionen (1)

„Data mining ist definiert als der Prozess der Entdeckung von Mustern in Daten. Der Prozess muss automatisch oder (meistens) semiautomatisch sein. Die entdeckten Muster müssen bedeutungsvoll sein, da sie zu einem Vorteil, in der Regel ökonomischer Natur, führen. Die Daten sind unvermeidlich in erheblichen Mengen vorhanden.“

... und den des (maschinellen) Lernens:

„Dinge lernen, wenn sie ihr Verhalten auf eine Weise ändern, die es ihnen ermöglicht, in der Zukunft besser zu performen.“

Witten et al. (2013). *Data Mining*. Morgan Kaufmann.

Begriffsdefinitionen (2)

„Insbesondere definieren wir maschinelles Lernen als eine Reihe von Methoden, die automatisch Muster in Daten erkennen und diese Muster dann verwenden, um zukünftige Daten zu vorhersagen oder um andere Arten von Entscheidungsfindungen unter Unsicherheit (...) zu betreiben.“

Murphy, K. (2012) *Machine Learning - A Probabilistic Perspective*. MIT-Press. (S. 1)

Begriffsdefinitionen (3)

„Das Problem der Suche nach Mustern in Daten ist ein grundlegendes und hat eine lange und erfolgreiche Geschichte. Zum Beispiel erlaubten die umfangreichen astronomischen Beobachtungen von Tycho Brahe im 16^{ten} Jahrhundert Johannes Kepler die Entdeckung der empirischen Gesetze der planetarischen Bewegung, was wiederum die Entwicklung der klassischen Mechanik vorantrieb. (...) Das Feld der Mustererkennung befasst sich mit der automatisierten Entdeckung von Regelmäßigkeiten, um Maßnahmen wie die Klassifizierung von Daten in verschiedene Kategorien zu ergreifen.“

Bishop, C. (2007) *Pattern Recognition and Machine Learning*. MIT-Press. Springer.

Wozu maschinelles Lernen?

- Das indizierte WWW enthält am 8.9.2018 mindestens 4.42 Milliarden Webseiten (www.worldwidewebsize.com)
- Die Sentinel-Satelliten der Europäischen Weltraumorganisation (ESA) senden im Jahr 2016 täglich 10 Terabyte an die Bodenstation in Oberpfaffenhofen/DE. (www.dlr.de, 7.9.2018)
- ...

Arten des maschinellen Lernens

- *Überwachtes Lernen (supervised learning)*

Ziel: Aufgrund einer Menge von Input-Output-Paaren (Trainingsdaten) eine Abbildung finden, die einem Input-Vektor \mathbf{x} jeweils einen Output y zuordnet.

- Klassifikation (y ist diskret)
- Regression (y ist kontinuierlich)

- *Unüberwachtes Lernen (unsupervised learning)*

Ziel: „Interessante Muster“ in den Daten zu finden

- Clustering
- Dimensionsreduktion

- *Bestärkendes Lernen (reinforcement learning)*

Ziel: Mit welcher Strategie kann ein virtueller Agent in einer dynamischen Umgebung eine vorgegebene Belohnungsfunktion maximieren?

- genetische Algorithmen
- dynamische Programmierung

Diagnose

In den 80er Jahren des letzten Jahrhunderts zeigte eine Untersuchung, dass Data Mining-Regeln zur Diagnose von Krankheiten bei Sojabohnenpflanzen eine bessere Erkennungsrate hatten als Regeln, die von Experten aufgestellt wurden.

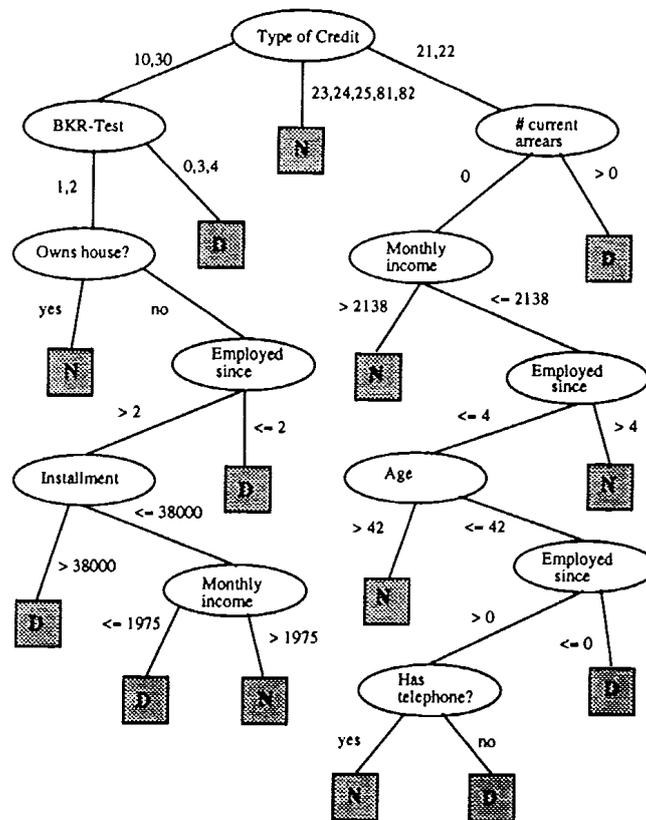


vom Diaporthe-Pilz befallene Sojabohne (*stem canker*)

<http://bulletin.ipm.illinois.edu/pastpest/articles/200218d.html> (31.12.2015)

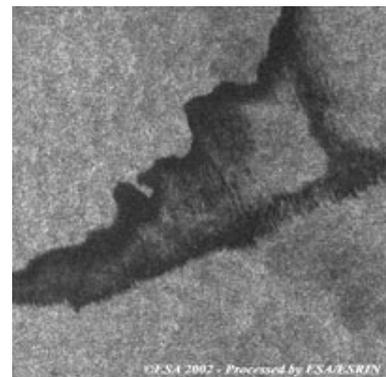
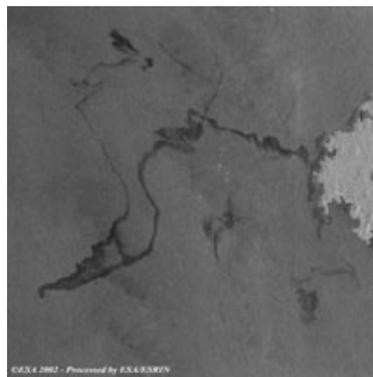
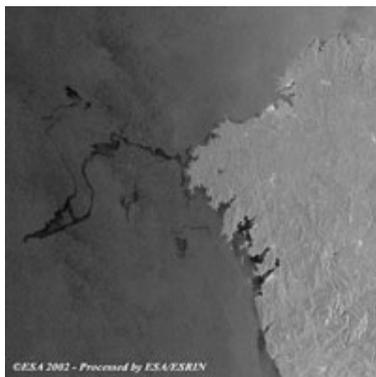
Kreditvergabe

In einer Fallstudie für Privatkredite (Feelder et al., 1995) wurden historische Kundendaten einer niederländischen Bank verwendet, um vorherzusagen, ob ein Kreditnehmer seinen Raten zurückzahlen wird. (D=Defaulter=säumiger Schuldner, N=Non Defaulter)



Analyse von Satellitenbildern

2002 bricht der mit 77 000 Tonnen Schweröl beladene Tanker „Prestige“ vor der Nordwestküste Spaniens entzwei und verursacht eine Naturkatastrophe.



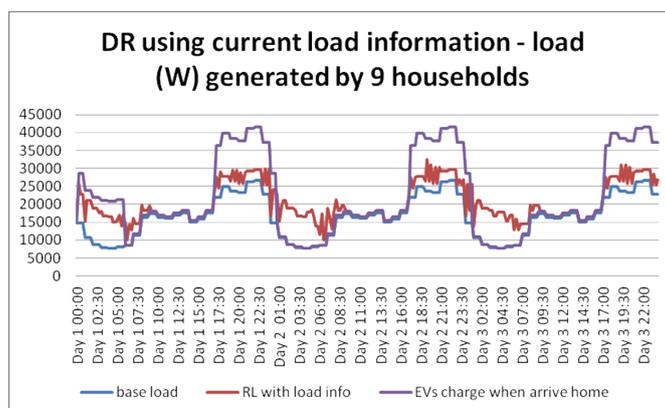
https://earth.esa.int/web/earth-watching/natural-disasters/oil-slicks/content/-/asset_publisher (30.12.2015)

Die schwarzen Flecken an der Küste im Bild links sind natürlichen Ursprungs.

Kann ein Computerprogramm den Unterschied zwischen „guten“ und „schlechten“ Flecken erkennen?

Optimierung von Angebot und Nachfrage in Stromnetzen

Durch „intelligente“ elektrische Geräte und Stromnetze lassen sich Schwankungen in der Stromnachfrage glätten. (DR = Demand Response, EV = Electric Vehicle)



<https://www.tcd.ie/futurecities/research/energy/residential-demand.php> (30.12.2015)

Marketing



<http://content.time.com/time/magazine/0,9263,7601110321,00.html> (30.12.2015)

Ethische Aspekte

Informationen über Geschlecht oder Rasse können für medizinische Fragen von Bedeutung sein. Wie verhält es sich jedoch bei der Vergabe eines Kredits oder bei der Auswahl eines neuen Mitarbeiters?

Selbst wenn solche Identifikationsmerkmale aus den Daten entfernt werden, könnte beispielsweise die Postleitzahl Auskunft über den sozialen Status einer Person geben. Darüber hinaus ist bekannt, dass 85% der Amerikaner durch ihr Geschlecht, die Postleitzahl ihres Wohnorts, und ihr vollständiges Geburtsdatum identifiziert werden können. (Witten, Frank & Hall, S. 33).