

Veri Madenciliği: Güncel Teknolojiler ve Araçlar

Veri madenciliği, büyük veri kümelerinden anlamlı bilgi, desen veya örüntüleri keşfetmeyi amaçlayan disiplinler arası bir alandır. Bu sunum, veri madenciliğinde kullanılan güncel teknolojiler ve araçları, web madenciliğini ve gelecek trendlerini ele alacaktır.

Veri madenciliği, istatistik, makine öğrenmesi, veri tabanı sistemleri ve diğer alanların kesişiminde yer alır. Günümüzde veri hacminin hızla artmasıyla birlikte, bu alandaki teknolojiler ve araçlar da sürekli gelişmektedir.

Dr. Öğr. Üyesi Alper Talha KARADENİZ

Veri Madenciliğine Genel Bakış



Tanım

Veri madenciliği, büyük veri kümelerinden anlamlı bilgi, desen veya örüntüleri keşfetmeyi amaçlayan disiplinler arası bir alandır. İstatistik, makine öğrenmesi, veri tabanı sistemleri ve diğer alanların kesişiminde yer alır.



Güncel Teknolojiler

Veri madenciliğinde kullanılan güncel teknoloji ve araçlar, işlenmesi zorlaşan, karmaşık ve büyüklüğü hızla artan veri setlerinden etkili ve verimli şekilde sonuçlar çıkarmayı hedefler.



Disiplinler Arası Yaklaşım

Veri madenciliği, farklı disiplinlerin bir araya gelmesiyle oluşan bir alandır ve çeşitli sektörlerde karar verme süreçlerini desteklemek için kullanılır.



Neden Güncel Teknolojiler ve Araçlar Önemli?



Günümüzde veri setleri artık petabayt düzeyine ulaşabilmektedir. Sosyal medya platformları, IoT sensörleri ve diğer kaynaklardan gelen büyük hacimli veriler, geleneksel tek makine üzerinde çalışan yöntemlerle işlenemez hale gelmiştir. Bu nedenle, performans ve ölçeklenebilirlik sorunlarını çözen güncel teknolojiler kritik önem taşımaktadır.





Programlama Dilleri

Python

Geniş makine öğrenmesi ve veri bilimi ekosistemine sahiptir: NumPy, pandas, scikit-learn, TensorFlow, PyTorch vb. kütüphaneler ile veri madenciliği projelerinde en çok tercih edilen dildir.

R

İstatistiksel analiz, veri görselleştirme ve çok sayıda paket (örn. caret, dplyr, ggplot2) içerir. Özellikle istatistiksel veri madenciliği uygulamalarında güçlü bir araçtır.

Java

Büyük ölçekli kurumsal uygulamalar ve dağıtık sistemlerle entegrasyonda popülerdir. Hadoop ekosistemi ile uyumlu çalışması nedeniyle büyük veri projelerinde tercih edilir.

Scala

Spark ekosistemi içinde (Spark Core, Spark MLlib) kullanılan, fonksiyonel ve nesne yönelimli dilleri birleştiren teknolojidir. Dağıtık veri işleme için optimize edilmiştir.

Kütüphaneler ve Framework'ler



scikit-learn

Makine öğrenmesi ve veri madenciliği için temel algoritmaları (sınıflandırma, regresyon, kümeleme vb.) içerir. Python ekosisteminin en popüler ML kütüphanesidir.



Spark MLlib

Apache Spark
çerçevesi üzerinde
büyük veriyi dağıtık
olarak işlemek için
makine öğrenmesi
kütüphanesi. Yüksek
performanslı dağıtık
hesaplama sağlar.



TensorFlow & PyTorch

Derin öğrenme (Deep Learning) modellerini geliştirmek için kullanılır. Karmaşık sinir ağları oluşturmak ve eğitmek için güçlü araçlar sunarlar.



Keras

TensorFlow üzerine inşa edilen yüksek seviyeli API. Derin öğrenme projelerinde hızlı prototipleme ve kullanım kolaylığı sunar.





4

3

Dağıtık ve Büyük Veri Teknolojileri

Hadoop Ekosistemi

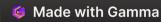
Temel bileşenler: HDFS (Hadoop Distributed File System), YARN, MapReduce. Büyük veri setlerini parçalara bölerek çoklu düğümlerde saklama ve işleme olanağı sağlar.

NoSQL Veritabanları

MongoDB, Cassandra, HBase gibi dağıtık ve yatayda ölçeklenebilen veritabanları. Yarı yapılandırılmış veya yapılandırılmamış verilerde veri madenciliği için kullanılabilir.

Apache Spark

Hadoop MapReduce'a göre daha hızlı bir bellek içi (inmemory) işlem modeli sunar. Spark SQL, Spark Streaming, MLlib, GraphX gibi bileşenleri vardır.



Bulut Tabanlı Çözümler

AWS Sagemaker

Gömülü veri madenciliği, otomatik model eğitimi, dağıtık işlem. S3 ve Redshift gibi depolama ve veri ambarı hizmetleriyle entegre çalışır.



Microsoft Azure Machine Learning

Sürükle-bırak arayüzlü veri madenciliği modülleri içerir. Azure bulut altyapısıyla ölçeklenebilir yapıda.

Google Cloud AI Platform

TensorFlow ve diğer ML araçlarını entegre eden bulut platformu. Büyük veri analitiği için BigQuery ve Dataflow servisleriyle etkileşim.

Bulut tabanlı çözümler, veri madenciliği projelerinde altyapı yönetimi ihtiyacını ortadan kaldırarak, araştırmacıların ve veri bilimcilerin daha hızlı sonuç almalarını sağlar. Ayrıca, ihtiyaca göre ölçeklenebilir kaynaklar sunarak maliyet optimizasyonu da sağlarlar.



Veri Madenciliğinde Kullanılan Hazır Programlar



Weka

Java tabanlı, akademik çevrelerde yaygın. Kapsamlı makine öğrenmesi algoritmaları, veri ön işleme araçları, görsel arayüz (GUI) içerir.



RapidMiner

Sürükle-birak şeklinde işlem akışları tasarlama (ETL, model eğitimi, validasyon). Kod yazmaya gerek duymadan birçok algoritmayı deneyebilirsiniz.



KNIME

Veri bilimi işlemlerini görsel iş akışlarıyla tasarlamayı kolaylaştıran açık kaynaklı araç. İstatistik, makine öğrenmesi, metin madenciliği, veri temizleme modülleri mevcuttur.

Bu hazır programlar, programlama bilgisi olmayan kullanıcıların da veri madenciliği projelerini gerçekleştirebilmesini sağlar. Görsel arayüzleri sayesinde karmaşık veri işleme ve analiz süreçleri daha anlaşılır hale gelir.

Web Madenciliği (Web Mining)



Web madenciliği, veri madenciliği tekniklerini web kaynakları (örn. web sayfaları, sunucu logları, sosyal medya içerikleri) üzerinde uygulayarak anlamlı bilgi çıkarmayı amaçlar. Bu alan, e-ticaret sitelerinin kullanıcı davranışlarını analiz etmek, arama motorlarının sonuçlarını iyileştirmek ve sosyal ağ analizleri yapmak için yaygın olarak kullanılmaktadır.



Gelecek Trendleri ve Sonuç

100+

60%

5x

Veri Kaynakları

Veri çeşitliliği artıyor: metin, video, sensör, sosyal medya, IoT

Otomasyon

Veri madenciliği süreçlerinde otomatikleşme oranı

Hız Artışı

Güncel teknolojilerle veri işleme hızındaki artış

Veri madenciliği, güncel teknolojiler ve araçlar sayesinde daha erişilebilir ve daha güçlü hâle geldi. Artan veri hacmi ve veri çeşitliliği (metin, video, sensör, sosyal medya, IoT vb.) nedeniyle, veri madenciliği alanında büyük veri altyapıları, derin öğrenme ve bulut çözümleri giderek önem kazanıyor.

Gelecekte, otomatik makine öğrenmesi (AutoML), kenar bilişim (edge computing) ve federe öğrenme gibi teknolojilerin veri madenciliği alanında daha fazla yer bulması beklenmektedir. Bu gelişmeler, veri bilimcilerin daha karmaşık problemleri daha hızlı çözebilmesini sağlayacaktır.