Big Data Analytics 6hp

http://www.ida.liu.se/~patla00/courses/BDA

Teachers

Lectures: Patrick Lambrix, Christoph Kessler, Jose Pena, Valentina Ivanova, Labs: Zlatan Dragisic, Huanyu Li NSC: Rickard Armiento

Course literature

- Articles (on web)
- Lab descriptions (on web)

Data and Data Storage

Data and Data Storage

Database / Data source

- One (of several) ways to store data in electronic format
- Used in everyday life: bank, hotel reservations, library search, shopping

Databases / Data sourcces

- Database management system (DBMS): a collection of programs to create and maintain a database
- Database system = database + DBMS

Databases / Data sources



What information is stored?

- Model the information
 - Entity-Relationship model (ER)
 - Unified Modeling Language (UML)

What information is stored? - ER

- entities and attributes
- entity types
- key attributes
- relationships
- cardinality constraints

EER: sub-types

1 tgctacccgc gcccgggctt ctggggtgtt ccccaaccac ggcccagccc tgccacaccc 61 cccgcccccg gcctccgcag ctcggcatgg gcgcgggggt gctcgtcctg ggcgcctccg 181 tgctggtgcc cgcgtcgccg cccgcctcgt tgctgcctcc cgccagcgaa agccccgagc 241 cgctgtctca gcagtggaca gcgggcatgg gtctgctgat ggcgctcatc gtgctgctca 301 tcgtggcggg caatgtgctg gtgatcgtgg ccatcgccaa gacgccgcgg ctgcagacgc 361 tcaccaacct cttcatcatg tccctggcca gcgccgacct ggtcatgggg ctgctggtgg 421 tgccgttcgg ggccaccatc gtggtgtggg gccgctggga gtacggctcc ttcttctgcg 481 agetgtggac etcagtggac gtgetgtgeg tgacggecag categagace etgtgtgtca 541 ttgccctgga ccgctacctc gccatcacct cgcccttccg ctaccagagc ctgctgacgc 601 gcgcgcgggc gcggggcctc gtgtgcaccg tgtgggccat ctcggccctg gtgtccttcc 661 tgcccatcct catgcactgg tggcgggcgg agagcgacga ggcgcgccgc tgctacaacg 721 accccaagtg ctgcgacttc gtcaccaacc gggcctacgc catcgcctcg tccgtagtct 781 ccttctacgt gcccctgtgc atcatggcct tcgtgtacct gcgggtgttc cgcgaggccc 841 agaagcaggt gaagaagatc gacagctgcg agcgccgttt cctcggcggc ccagcgcggc 901 cgccctcgcc ctcgccctcg cccgtccccg cgccgcgcc gccgcccgga ccccgcgcc 961 ccgccgccgc cgccgccacc gcccgctgg ccaacgggcg tgcgggtaag cggcggccct 1021 cgcgcctcgt ggccctacgc gagcagaagg cgctcaagac gctgggcatc atcatgggcg 1081 tcttcacgct ctgctggctg cccttcttcc tggccaacgt ggtgaaggcc ttccaccgcg 1141 agctggtgcc cgaccgcctc ttcgtcttct tcaactggct gggctacgcc aactcggcct 1201 tcaaccccat catctactgc cgcagccccg acttccgcaa ggccttccag ggactgctct 1261 gctgcgcgcg cagggctgcc cgccggcgcc acgcgaccca cggagaccgg ccgcgcgcct 1321 cgggctgtct ggcccggccc ggacccccgc catcgcccgg ggccgcctcg gacgacgacg 1381 acgacgatgt cgtcggggcc acgccgcccg cgcgcctgct ggagccctgg gccggctgca 1441 acggcggggc ggcggcggac agcgactcga gcctggacga gccgtgccgc cccggcttcg 1501 cctcggaatc caaggtgtag ggcccggcgc ggggcgcgga ctccgggcac ggcttcccag 1561 gggaacgagg agatctgtgt ttacttaaga ccgatagcag gtgaactcga agcccacaat 1621 cctcgtctga atcatccgag gcaaagagaa aagccacgga ccgttgcaca aaaaggaaag 1681 tttgggaagg gatgggagag tggcttgctg atgttccttg ttg

DEFINITION ACCESSION SOURCE ORGANISM human REFERENCE AUTHORS

TITLE

REFERENCE AUTHORS TITLE

Homo sapiens adrenergic, beta-1-, receptor NM_000684

1

Frielle, Collins, Daniel, Caron, Lefkowitz, Kobilka

Cloning of the cDNA for the human beta 1-adrenergic receptor

2

Frielle, Kobilka, Lefkowitz, Caron Human beta 1- and beta 2-adrenergic receptors: structurally and functionally related receptors derived from distinct genes

Entity-relationship



Databases / Data sources



How is the information stored? (high level) How is the information accessed? (user level)



IR - formal characterization

Information retrieval model: (D,Q,F,R)

- D is a set of document representations
- Q is a set of queries
- F is a framework for modeling document representations, queries and their relationships
- R associates a real number to documentquery-pairs (ranking)

IR - Boolean model

	adrenergic	cloning	receptor	r	
Doc1	yes	yes	no	>	(1 1 0)
Doc2	no	yes	no	>	(0 1 0)

Q1: cloning and (adrenergic or receptor) --> (1 1 0) or (1 1 1) or (0 1 1) Result: Doc1 Q2: cloning and not adrenergic --> (0 1 0) or (0 1 1) Result: Doc2

IR - Vector model (simplified)



Semi-structured data



Semi-structured data - Queries

select source
from PROTEINDB.protein P
where P.accession = "NM 000684";

Relational databases

PROTEIN				REFERENCE	
PROTEIN-ID	ACCESSION	DEFINITION	SOURCE	PROTEIN-ID	ARTICLE-ID
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human	1 1	1 2

ARTICLE-AUTHOR		ARTICLE-TITLE	
ARTICLE-ID	AUTHOR	ARTICLE-ID	TITLE
$ \begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 2 \\ 2 \\ 2 \\ 2 \\ 2 \\ 2 \end{array} $	Frielle Collins Daniel Caron Lefkowitz Kobilka Frielle Kobilka Lefkowitz Caron	1 2	Cloning of the cDNA for the human beta 1-adrenergic receptor Human beta 1- and beta 2- adrenergic receptors: structurally and functionally related receptors derived from distinct genes

Relational databases - SQL

select source
from protein
where accession = NM_000684;

PROTEIN

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human

Evolution of Database Technology

- 1960s:
 - Data collection, database creation, IMS and network DBMS
- 1970s:
 - □ Relational data model, relational DBMS implementation
- 1980s:
 - □ Advanced data models (extended-relational, OO, deductive, etc.)
 - Application-oriented DBMS (spatial, temporal, multimedia, etc.)
- 1990s:
 - Data mining, data warehousing, multimedia databases, and Web databases
- 2000s
 - Stream data management and mining
 - Data mining and its applications
 - Web technology (XML, data integration) and global information systems
 - NoSQL databases

Knowledge bases

- (F) source(NM_000684, Human)
 (R) source(P?,Human) => source(P?,Mammal)
 (R) source(P?,Mammal) => source(P?,Vertebrate)
- Q: ?- source(NM_000684, Vertebrate) A: yes
- Q: ?- source(x?, Mammal) A: x? = NM_000684

Interested in more?

- 732A57 Database Technology (relational databases)
- TDDD43 Advanced data models and databases
 - (IR, semi-structured data, DB, KB)

732A47 Text mining (includes IR)

Analytics

Analytics

Discovery, interpretation and communication of meaningful patterns in data

Analytics - IBM

What is happening? Descriptive Discovery and explanation Why did it happen? Diagnostic Reporting, analysis, content analytics What could happen? Predictive Predictive analytics and modeling What action should I take? Prescriptive **Decision management** What did I learn, what is best? Cognitive

Analytics - Oracle

- Classification
- Regression
- Clustering
- Attribute importance
- Anomaly detection
- Feature extraction and creation
- Market basket analysis

Why Analytics?

. . .

- The Explosive Growth of Data
 - Data collection and data availability
 - Automated data collection tools, database systems, Web, computerized society
 - □ Major sources of abundant data
 - Business: Web, e-commerce, transactions, stocks, ...
 - Science: Remote sensing, bioinformatics, scientific simulation,
 - Society and everyone: news, digital cameras, YouTube
- We are drowning in data, but starving for knowledge!

Ex.: Market Analysis and Management

- Where does the data come from?—Credit card transactions, loyalty cards, discount coupons, customer complaint calls, plus (public) lifestyle studies
- Target marketing
 - □ Find clusters of "model" customers who share the same characteristics: interest, income level, spending habits, etc.
 - Determine customer purchasing patterns over time
- Cross-market analysis—Find associations/co-relations between product sales, & predict based on such association
- Customer profiling—What types of customers buy what products (clustering or classification)
- Customer requirement analysis
 - Identify the best products for different groups of customers
 - Predict what factors will attract new customers
- Provision of summary information
 - Multidimensional summary reports
 - □ Statistical summary information (data central tendency and variation)

Ex.: Fraud Detection & Mining Unusual Patterns

- Approaches: Clustering & model construction for frauds, outlier analysis
- Applications: Health care, retail, credit card service, telecomm.
 - □ <u>Auto insurance</u>: ring of collisions
 - □ <u>Money laundering</u>: suspicious monetary transactions
 - Medical insurance
 - Professional patients, ring of doctors, and ring of references
 - Unnecessary or correlated screening tests
 - Telecommunications: phone-call fraud
 - Phone call model: destination of the call, duration, time of day or week. Analyze patterns that deviate from an expected norm
 - Anti-terrorism



Data Mining: Classification Schemes

General functionality

Descriptive data mining

Predictive data mining

- Concept/class description:
 - Characterization: summarizing the data of the class under study in general terms
 - E.g. Characteristics of customers spending more than 10000 sek per year
 - Discrimination: comparing target class with other (contrasting) classes
 - E.g. Compare the characteristics of products that had a sales increase to products that had a sales decrease last year

- Frequent patterns, association, correlations
 - □ Frequent itemset
 - Frequent sequential pattern
 - Frequent structured pattern
 - □ E.g. buy(X, "Diaper") → buy(X, "Beer") [support=0.5%, confidence=75%]
 confidence: if X buys a diaper, then there is 75% chance that X buys beer
 support: of all transactions under consideration 0.5% showed that diaper and
 beer were bought together
 - □ E.g. Age(X, "20..29") and income(X, "20k..29k") → buys(X, "cd-player")
 [support=2%, confidence=60%]

- Classification and prediction
 - Construct models (functions) that describe and distinguish classes or concepts for future prediction.

The derived model is based on analyzing training data

- data whose class labels are known.

 E.g., classify countries based on (climate), or classify cars based on (gas mileage)

Predict some unknown or missing numerical values

Cluster analysis

Class label is unknown: Group data to form new classes, e.g., cluster customers to find target groups for marketing

- □ Maximizing intra-class similarity & minimizing interclass similarity
- Outlier analysis
 - Outlier: Data object that does not comply with the general behavior of the data
 - □ Noise or exception? Useful in fraud detection, rare events analysis
- Trend and evolution analysis
 - Trend and deviation

Interested in more?

732A95 Introduction to machine learning
 TDDD41 Data mining – clustering and association analysis

Big Data

Big Data

So large data that it becomes difficult to process it using a 'traditional' system

Big Data – 3Vs

Volumesize of the data

Volume - examples

- Facebook processes 500 TB per day
- Walmart handles 1 million customer transaction per hour
- Airbus generates 640 TB in one fligth (10 TB per 30 minutes)
- 72 hours of video uploaded to youtube every minute
- SMS, e-mail, internet, social media

What Happens in an Internet Minute?



https://y2socialcomputing.files.wordpress.com/2012/06/

social-media-visual-last-blog-post-what-happens-in-an-internet-minute-infographic.jpg

Big Data – 3Vs

- Volume
 - size of the data
- Variety
 - □ type and nature of the data
 - text, semi-structured data, databases, knowledge bases



Linking Open Data cloud diagram 2014, by Max Schmachtenberg, Christian Bizer, Anja Jentzsch and Richard Cyganiak. http://lod-cloud.net/

Linked open data of US government

Format (# Datasets)

- HTML (27005)
- XML (24077)
- PDF (19628)
- CSV (10058)
- JSON (8948)
- RDF (6153)
- JPG (5419)
- WMS (5019)
- Excel (3389)
- WFS (2781)

http://catalog.data.gov/

Big Data – 3Vs

- Volume
 - size of the data
- Variety
 - type and nature of the data
- Velocity
 - speed of generation and processing of data

Velocity - examples

- Traffic data
- Financial market
- Social networks



http://www.ibmbigdatahub.com/infographic/four-vs-big-data



Big Data – other Vs

Variability

inconsistency of the data

Veracity

quality of the data

Value

_ _ _

useful analysis results

BDA system architecture

Specialized services for domain A Specialized services for domain B

Big Data Services Layer

Knowledge Management Layer

Data Storage and Management Layer

BDA system architecture

- □ Large amounts of data, distributed environment
- Unstructured and semi-structured data
- Not necessarily a schema
- Heterogeneous
- Streams
- Varying quality

Data Storage and Management Layer

Data Storage and management – this course

Data storage: □ NoSQL databases □ OLTP vs OLAP Horizontal scalability Consistency, availability, partition tolerance Data management □Hadoop

Data management systems

BDA system architecture

Semantic technologies

- Integration
- Knowledge acquisition

Knowledge Management Layer

Knowledge management – this course

- Not a focus topic in this course
- For semantic and integration approaches see TDDD43

BDA system architecture

Analytics services for Big Data

Big Data Services Layer

Big Data Services – this course

Big data versions of analytics/data mining algorithms



Course overview

- Databases for Big Data (lectures + lab)
- Parallel algorithms for processing Big Data (lectures + lab)
- Machine Learning for Big Data (lectures + lab)

Visit to National Supercomputer Centre

Credits for the course

Written exam: May 10, 8-12
 LiU: sign up for 732A54 (ca april 20-30)
 Others: contact with supervisor

Labs

HARD DEADLINE: Labs approved by April 30. (No guarantee NSC resources available after April.)

Visit to NSC

Leave from here 16:00.

Or

Be in G34 latest 16:15.

My own interest and research

- Modeling of data
 - Ontologies
- Ontology engineering
 - Ontology alignment (Winner Anatomy track OA)
 - (Winner Anatomy track OAEI 2008 / Organizar OAEI tracks since 2012)
 - Organizer OAEI tracks since 2013)
 - Ontology debugging
 - (Founder and organizer WoDOOM/CoDeS 2012-2016)
 - Ontologies and databases for Big Data
- Former work: knowledge representation, data integration, knowledge-based information retrieval, object-centered databases

http://www.ida.liu.se/~patla00/research.shtml

https://www.youtube.com/watch?v=LrNIZ7-SMPk