

TADA

practicalities & more

on DM

24 April 2014



More on Data Mining as a Science

DM as method development

- Data mining develops methods for scientists
 - C.f. mathematics or statistics
- The research of DM in universities doesn't follow the scientific paradigm
 - But that doesn't make it a voodoo science
 - ...the applications of DM are another story

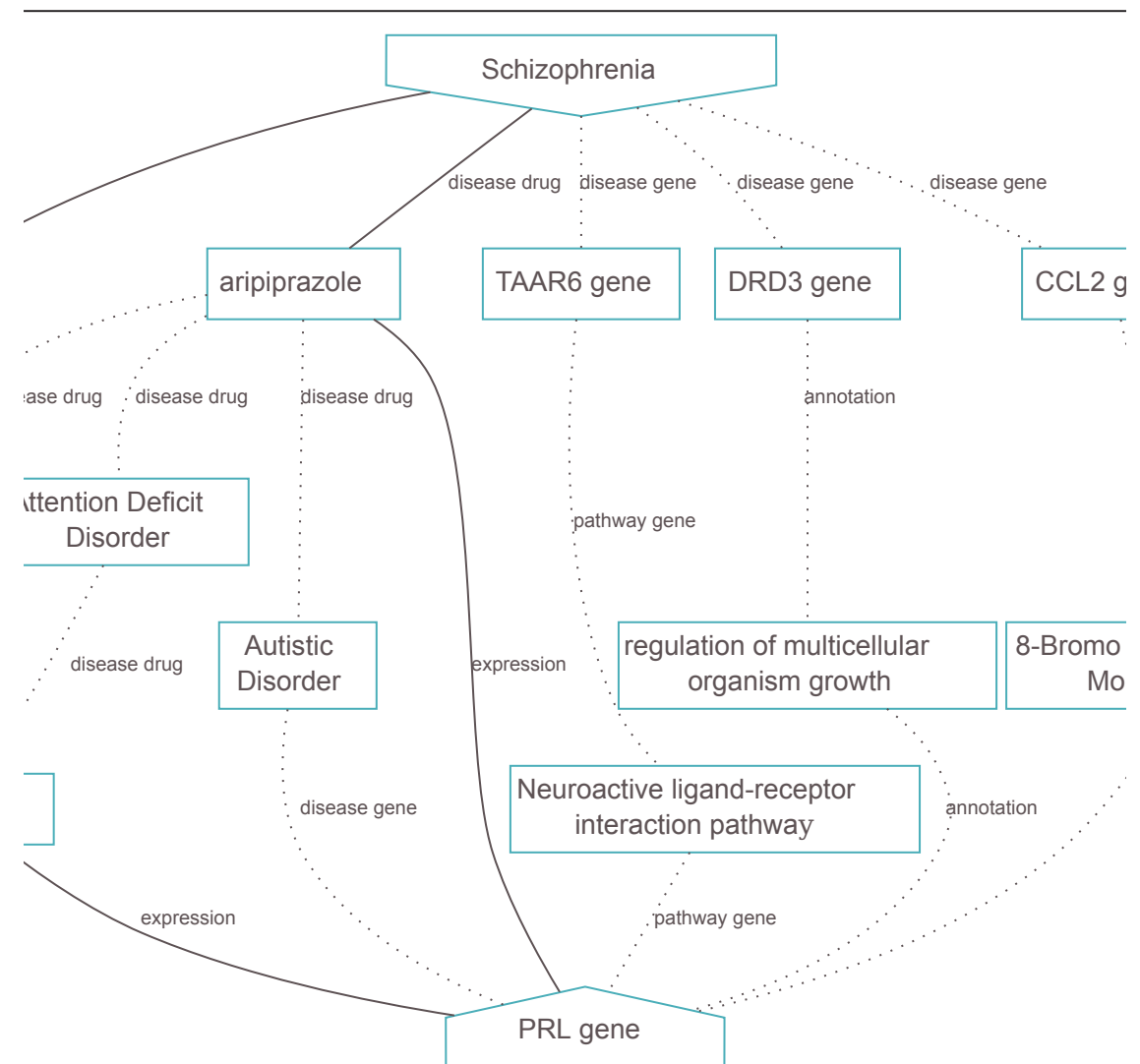
Of DM, ML, and Stat

- One trichotomy:
 - Statistics studies how reliable inferences can be drawn from imperfect data
 - ML develops technology of automated induction
 - DM is the art of extracting useful patterns from large bodies of data

Data Mining success stories

Bioinformatics

- BioGraph provides automated inference of functional hypotheses
- E.g. which genes are most potential to be associated with certain diseases



presentation of the top ten automatically generated hypotheses supporting the suspected association between Schizophrenia and the PRL gene. The solid and dotted line styles represent the importance of the link in descending order, that is, the most important links are solid lines. The nodes represent gene concepts while performing random walks from the source schizophrenia concept in the curated knowledge bases, annotated with their semantic meanings and enriched by their

Making money

- “Recommended for you”
- “Others often bought also”
- All of modern targeted advertisement is based on some type of data mining

Wird oft zusammen gekauft



Preis für alle drei: **EUR 3.218,65**

[Alle drei in den Einkaufswagen](#)

[Verfügbarkeit und Versanddetails anzeigen](#)

- ✓ **Dieser Artikel:** Canon EOS 5D Mark III SLR-Digitalkamera (22 Megapixel, CMOS-Sensor, 8,1 cm (3,2 Zoll) Display, ... **EUR 2.854,00**
- ✓ Canon EOS BG-E11 Batteriegriff für Canon EOS 5D Mark III und LP-E6 **EUR 299,90**
- ✓ Canon Kamera Akkupack LP-E6 (LI-ION) **EUR 64,75**

Hörer kauften auch:



Google AdWords

Obama's re-election

- Data of electorate was used to target the campaigning efforts where they count
- DM was also used to optimize fund-raising from small donations



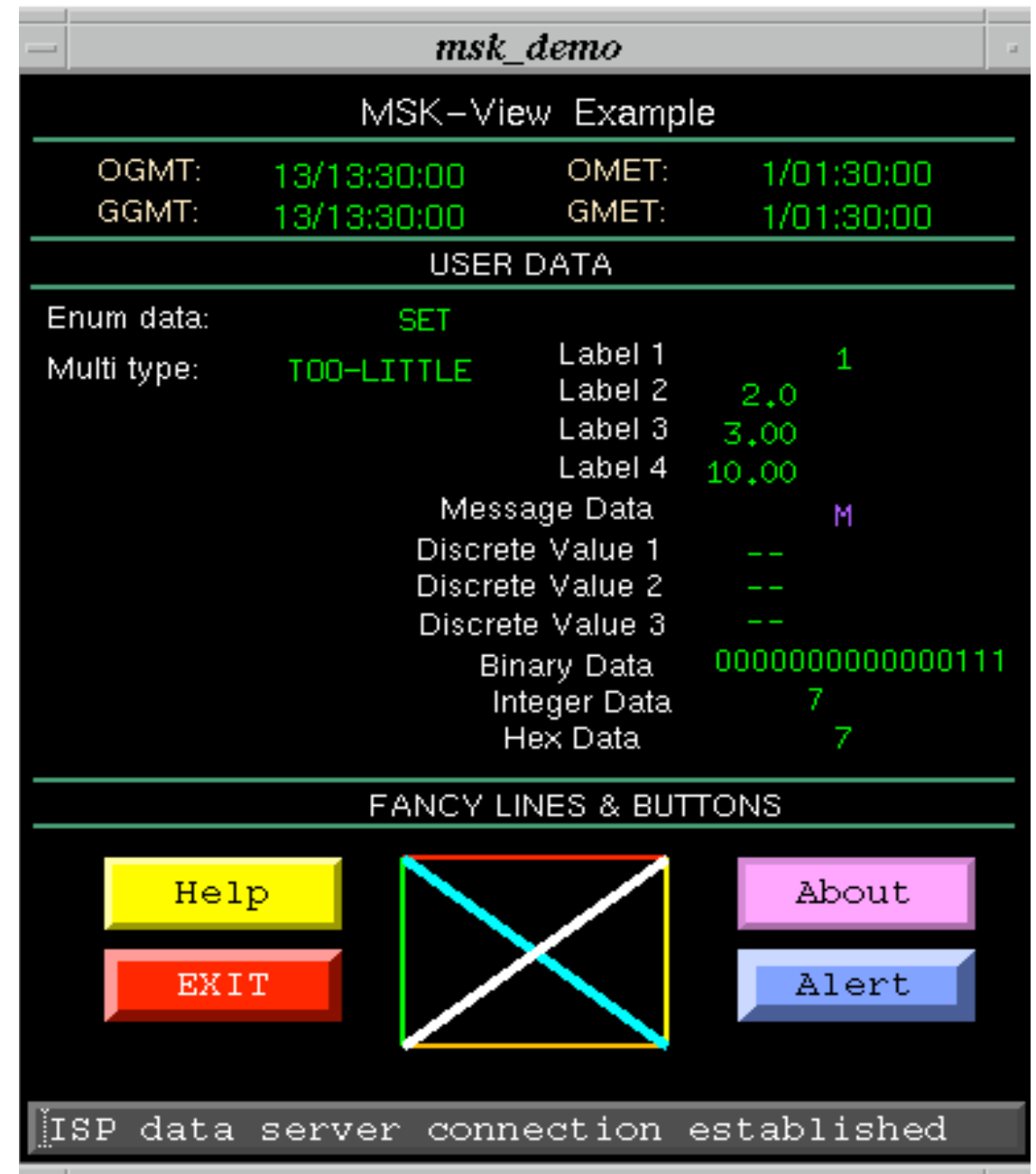
Church uses Big Data

- Evangelical Lutheran Church of Finland uses data mining to study its parishes
- What type of people live in which geographical areas?

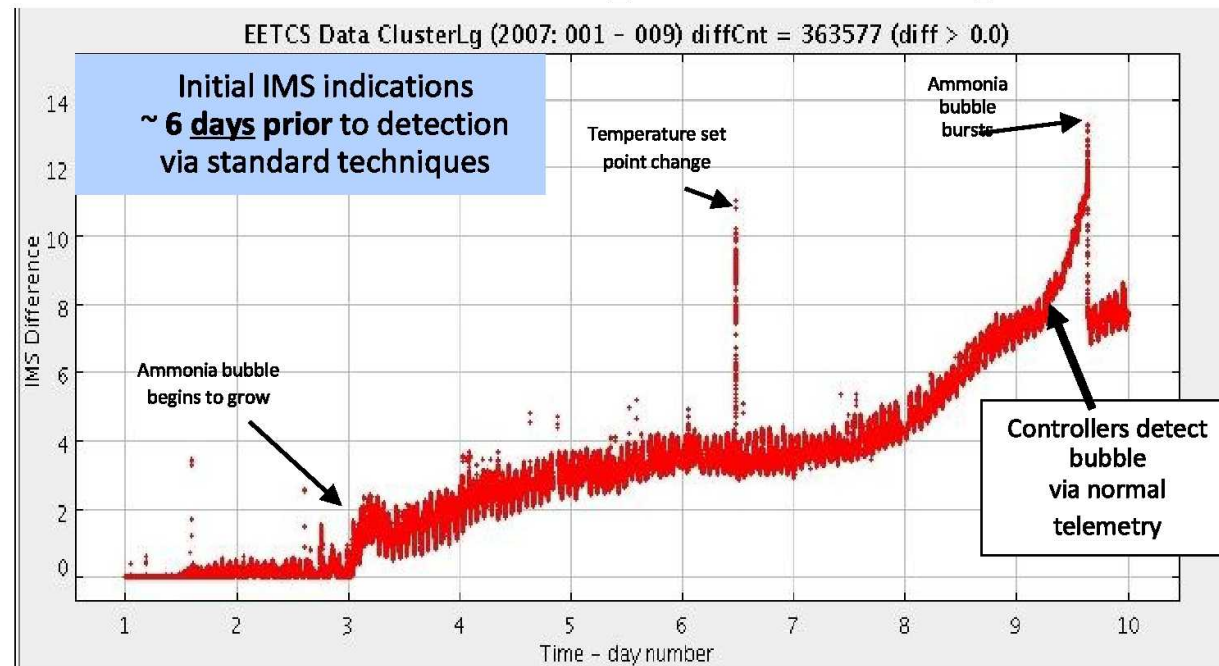


Space program safety

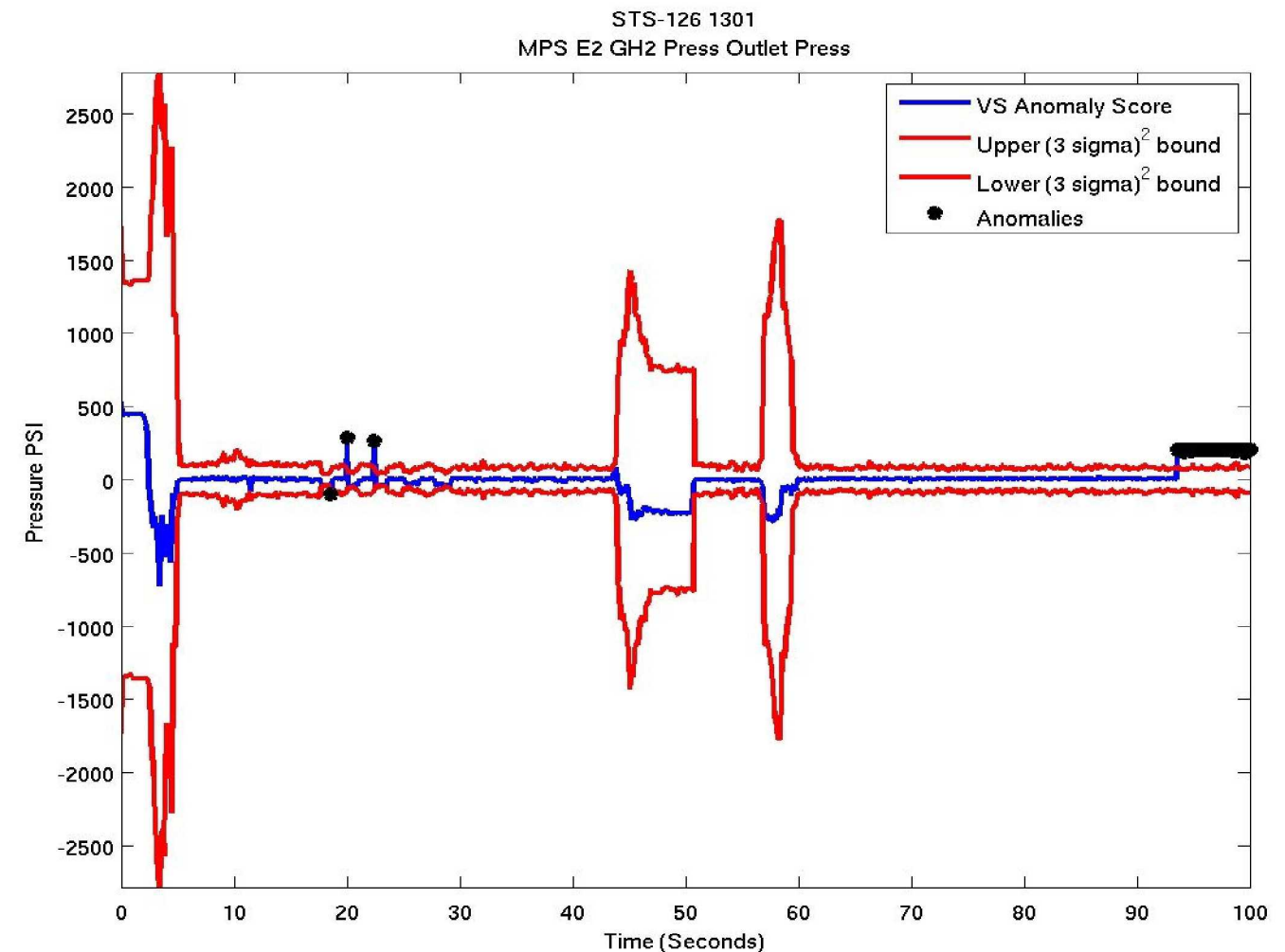
- ORCA searches outliers from sensor readings by comparing parameter-value vectors to their neighbors
- IMS builds a model of normal variance of sensor readings to detect anomalies



More on IMS



- In early January 2007, ISS Early External Thermal Control System developed an ammonia gas bubble
- Bubble noted by ISS controllers only ~9 hours before it “burst” and dissipated back into liquid



Practicalities

Schedule

Month	Day	Lecture topic	Assignments
April	17	Intro	
	24	Practicalities & where DM is used	1st assignment given out
May	1	<i>No lecture</i> (First of May)	
	8	Intro to Tensors	1st assignment DL, 2nd assignment given out
	15	Tensors in DM	
	22	Special topics in tensors	
	29	<i>No lecture</i> (Ascension day)	
June	5	MDL for pattern mining	2nd assignment DL, 3rd assignment given out
	12	Maximum entropy & iterative data mining	
	19	<i>No lecture</i> (Corpus Christi)	
	26	Kolmogorov complexity, cumulative entropy, and causality	
July	3	Graphs I	3rd assignment DL, 4th assignment given out
	10	Graphs II	
	17	Graphs III	
	24	Wrap-up	4th assignment DL

On Exam

- Day and place TBA
 - Most likely in early September
- Type TBA
- Final grade is based on the final exam and the assignments
 - Assignments also determine the eligibility to sit the final exam

On assignments: general

- 4 assignments
- Grading: fail, pass, excellent
- You can fail one assignment
 - 2 fails \Rightarrow course failed
- Every excellent gives $1/3$ point improvement on the final exam grade
 - But maximum of 1 full point (3 ex's)
- **You must pass the final exam to pass the course**

On assignments: requirements

- Assignments are to be written in proper academic-style English
- **Proper citations**
- You are given sources, but you can also use outside sources
 - Naturally must be mentioned
 - Plagiarism ⇒ failed assignment

On assignments: format

- Assignments need to be returned as PDF files by email
 - **No** .doc(x), .odt, .rtf, .txt, .xml, .html, .pages, .ps, .wp, or anything else
- No length limits — use the space you need
 - Probably most will need 3–4 pages...
- All PDFs must have name, matriculation number, email address, and clearly state the topic

On assignments: returning

- The assignments are returned by email to `tada14@mpi-inf.mpg.de`
- DL is 1600 hours on the stated day
 - **No** delays, no excuses, time based on the mail time stamp
- We'll acknowledge the submission that we receive before the lecture on the DL day

On assignments: grading

- Assignments are **not** for repeating what the papers say
 - We've read the papers already
- We expect you to discuss **and criticize** the sources, build connections, point out differences, provide new insights, etc.
- Some assignments are marked *hard*
 - This is taken into account when grading

First assignments

1. Did Tukey invent Data Mining?
2. (Don't) Believe the Hype
3. Big Data: The Best Thing since Sliced Bread or just Another Bottle of Snake Oil?
4. Where did the Candidates Go? (*Hard*)