

Winnen met Data Science

Joost N. Kok, Leiden Centre of Data Science

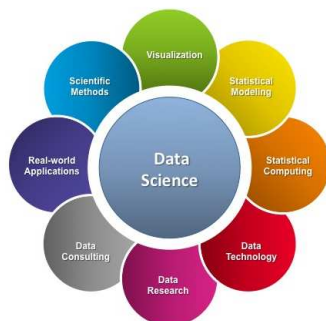


**Universiteit
Leiden**
The Netherlands

Discover the world at Leiden University

Overzicht

Data Science



Sport en Beweging

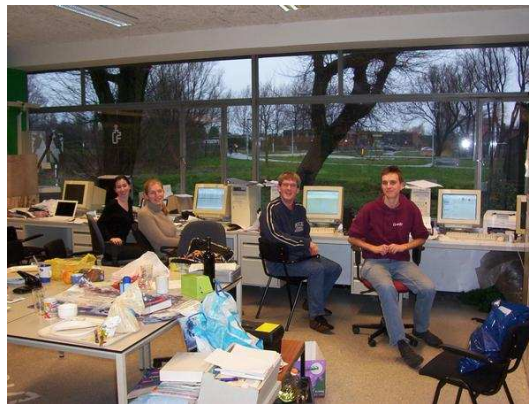


Discover the world at Leiden University



LIACS

- The Computer Science institute of Leiden University






Discover the world at Leiden University

LCDS Profile

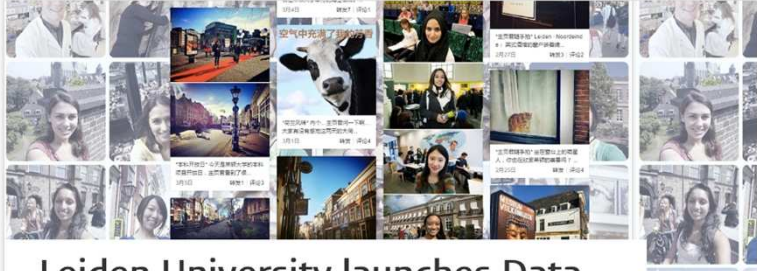
- **Data Science Dossier:**
<http://onderzoeksgebieden.leidenuniv.nl/data-science>
- Focus op
 - Fundamenteel
 - Standaarden
 - Wetenschappelijke Data

Discover the world at Leiden University


Universiteit Leiden
All

[Home](#)
[Research](#)
[Education](#)
[Academic staff](#)
[About us](#)
[Faculties](#)
[The Hague](#)
[Library](#)

Home > News > Leiden University launches Data Science research programme



Leiden University launches Data Science research programme

14 April 2016

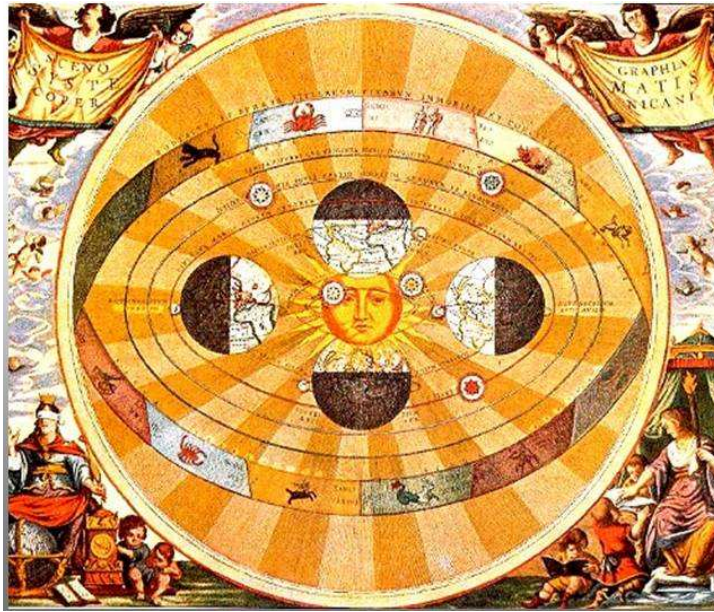
Leiden University is investing 4 million euros in a new Data Science research programme. This is a joint initiative of all the faculties, headed by Dean Geert de Snoo at the Faculty of Science. The programme will focus on Leiden scientific data.

Knowledge exchange is the essence

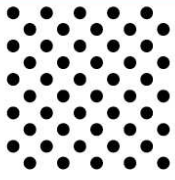
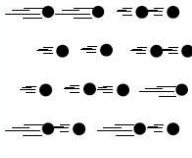
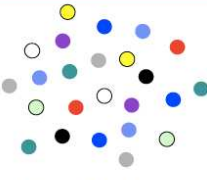

De Snoo: 'The Leiden Data Science programme is important for all our faculties, in particular for

Discover the world at Leiden University



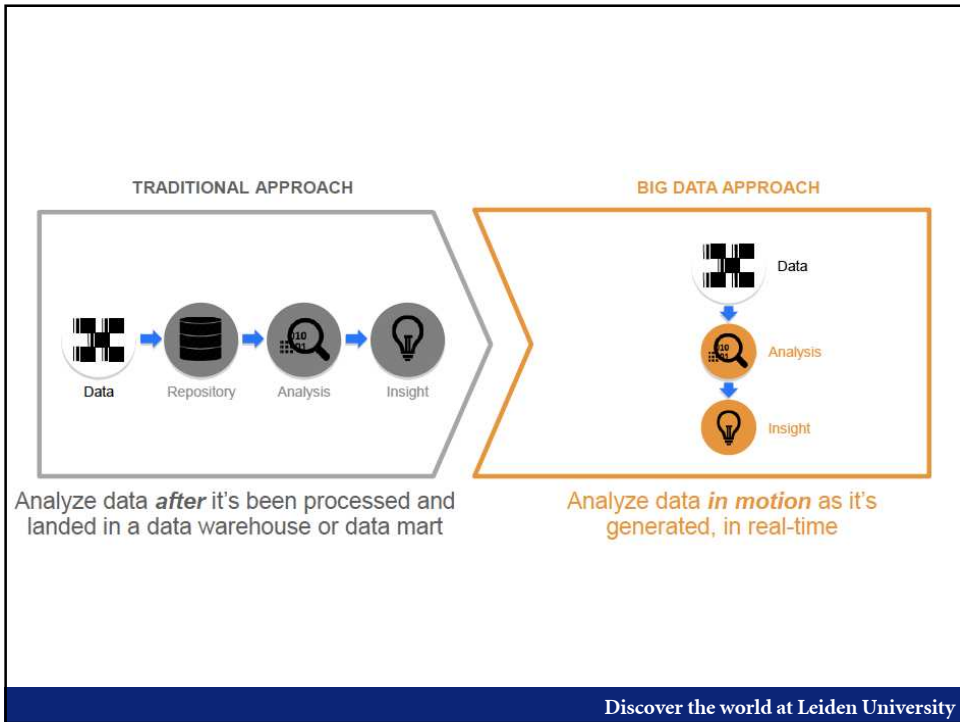
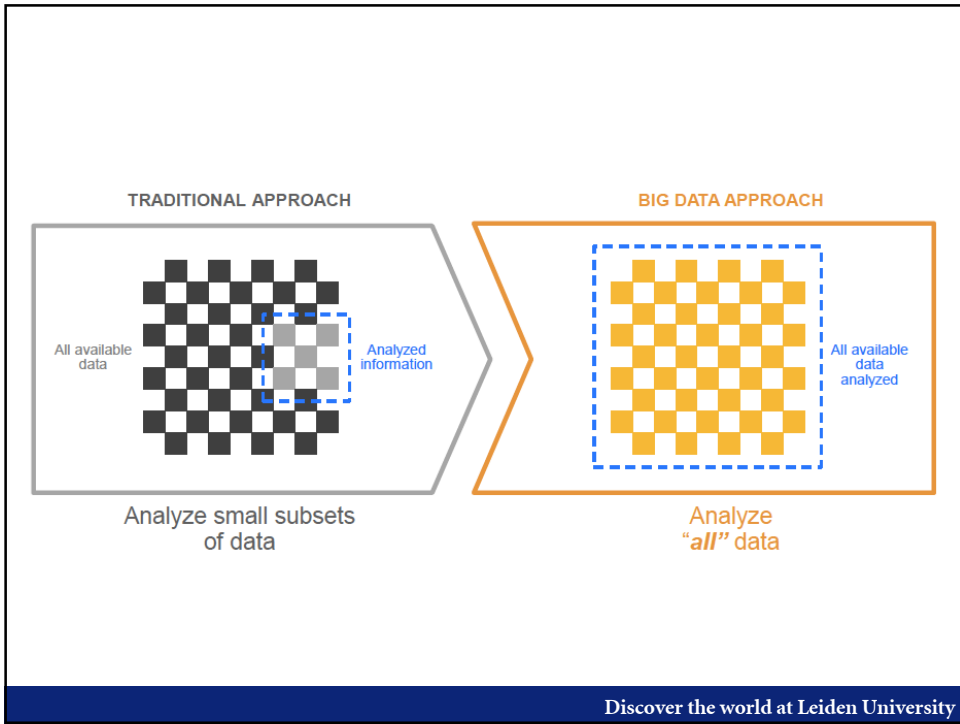


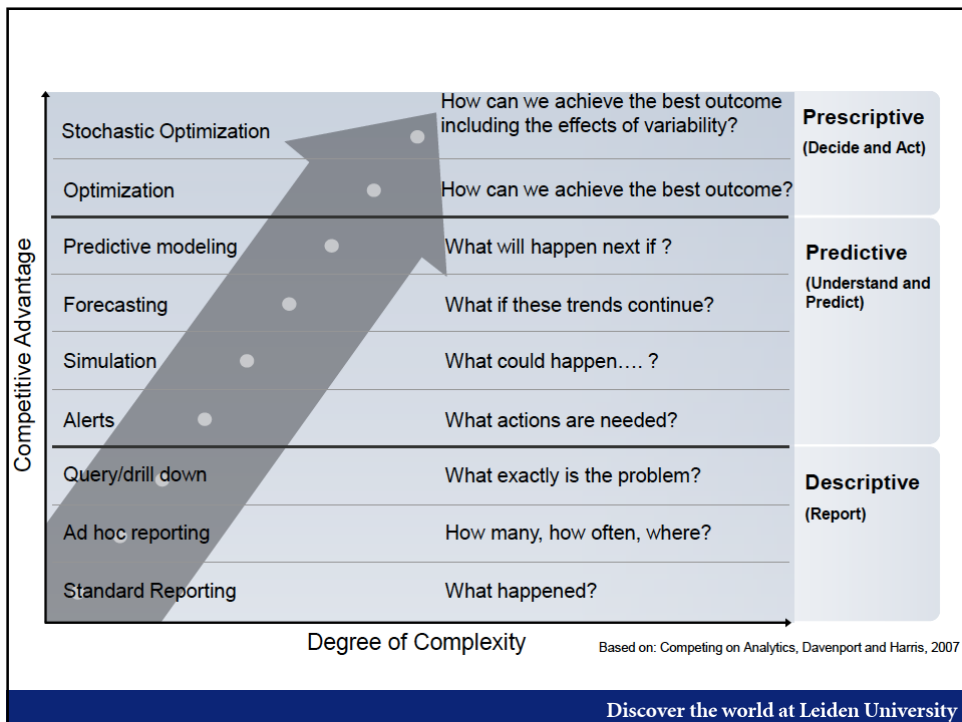
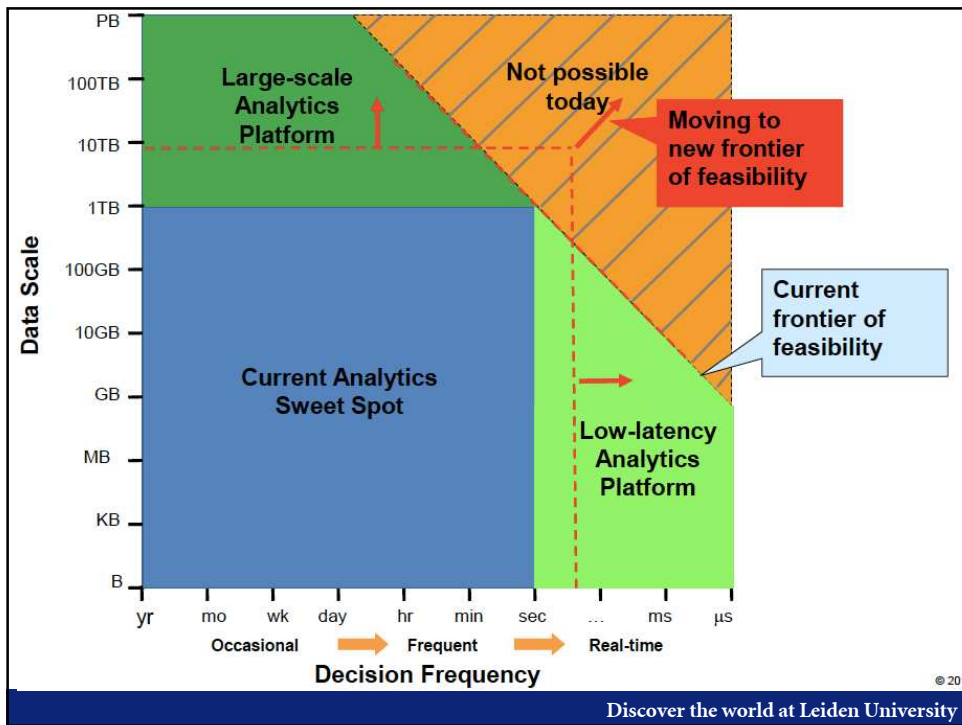
Discover the world at Leiden University

Volume	Velocity	Variety	Veracity
 <p>Data at Scale Terabytes to exabytes of existing data to process (e.g. CRM, ERP data, etc.)</p>	 <p>Data in Motion Streaming data, milliseconds to seconds to respond (e.g. data from smart sensors, mobile device, etc.)</p>	 <p>Data in Many Forms Structured, unstructured, text, multimedia (e.g. relational DB, images, free text, video, etc.)</p>	 <p>Data Uncertainty Uncertainty due to data error, inconsistency & incompleteness, ambiguities, model approximations (e.g. manual errors, device errors, models errors, etc.)</p>

Source: IBM GTO 2012

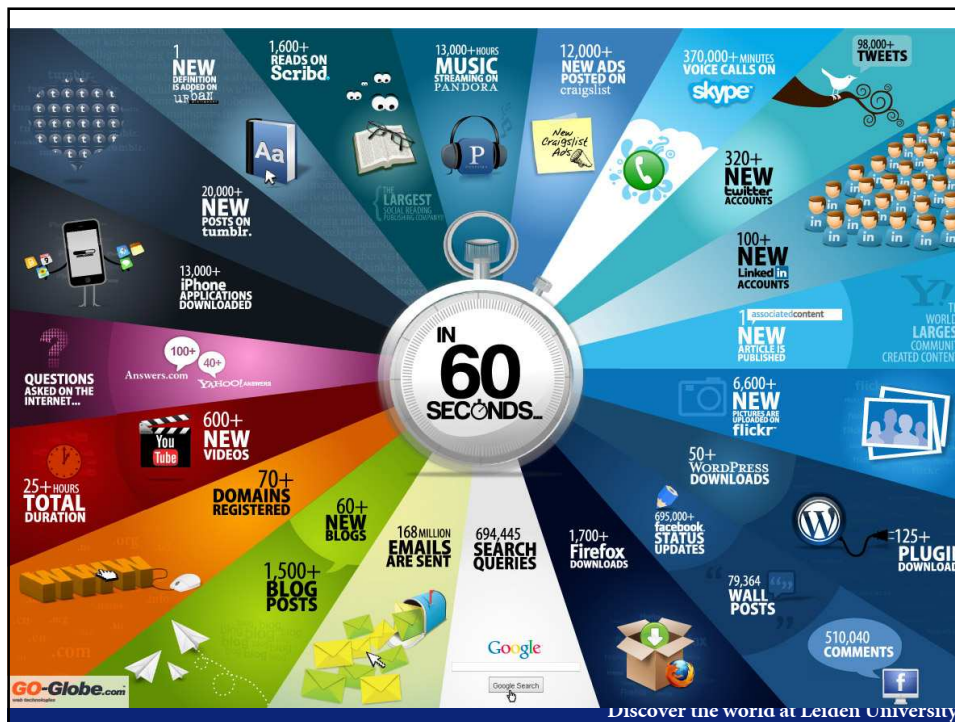
Discover the world at Leiden University





Data, data, data

- 90% of the world's data is created in the last two years
- 80% of the data is unstructured
- 1 trillion (10^{12}) connected devices generate 2.5 quintillion (10^{18}) bytes per day



Value of Data

- <http://tinyurl.com/valueofdata>

The screenshot shows a web form titled "What is your data worth?". At the top, there is a progress bar with five categories: DEMOGRAPHICS (active), FAMILY & HEALTH, PROPERTY, ACTIVITIES, and CONSUMER. Below the progress bar, there is a text box with the following text: "Data brokers scour public documents, such as birth records and motor vehicle reports, to compile basic data about individuals. It is likely they already know your:" followed by a list of checkboxes: Age, Gender, ZIP code, Employer, and Education level. Below this, there are several questions with radio buttons: "Are you a millionaire?" (No, Yes), "What is your job?" (a dropdown menu), "Are you engaged to be married?" (No, Yes), and "Are you?" (Recently married, Recently divorced, Empty nester). On the right side of the form, there is a large red box displaying "\$0.007" with the text "Current value of my data" below it. At the bottom right of the form, there is a "NEXT >" button.

Discover the world at Leiden University

Computational Turn

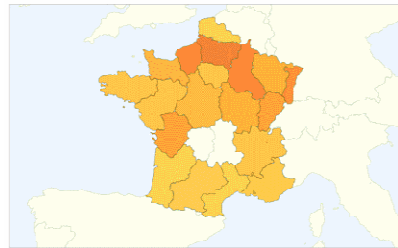
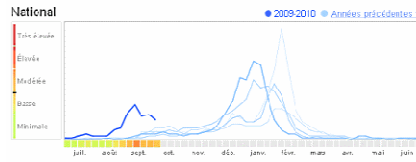
- Correlaties volgen direct uit de data
- (Heel) veel data beschikbaar
- Grote vooruitgang in methoden
- Redeneer over de correlaties om de oorzakelijke verbanden te vinden

Discover the world at Leiden University

Google: verspreiding griep

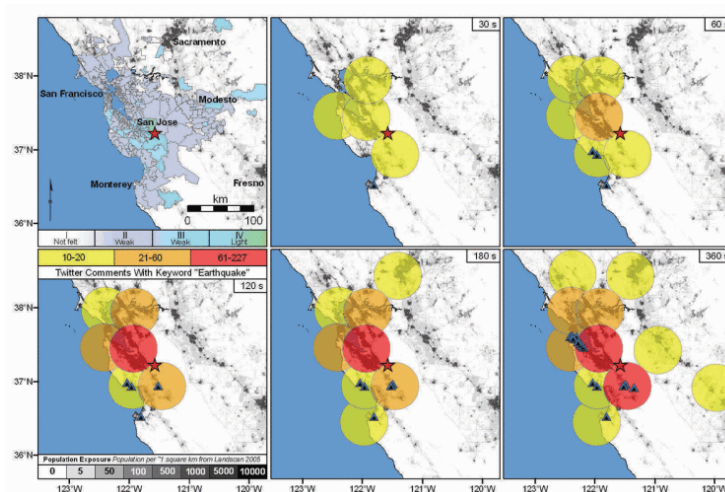
Évolution de la grippe - France

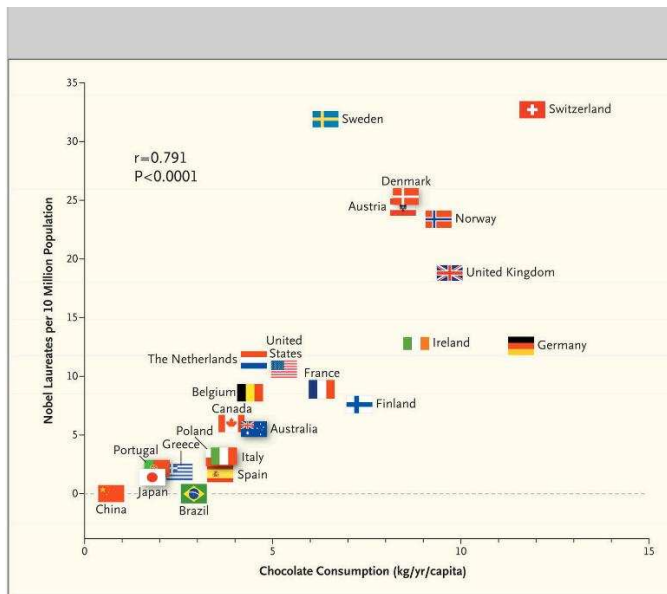
Certaines termes de recherche semblent être de bons indicateurs de la propagation de la grippe. Afin de vous fournir une estimation de la propagation du virus, ce site assemble donc des données relatives aux recherches lancées sur Google. [En savoir plus >](#)



Estimations faites réalisées à l'aide d'un modèle vérifié par rapport aux données officielles de propagation du virus. Données valides jusqu'au 8 octobre 2009.

Twitter: aardbeving





Data Science

- Data geschikt maken
 - Managen, labelen, combineren, ...
- Spelen met de data
 - Visualiseren, patroonherkenning, clusteren, ..
- Iets nuttigs mee doen
 - Classificeren, voorspellen, ...

MODERN DATA SCIENTIST


Data Scientist, the sexiest job of the 21st century, requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ☆ Machine learning
- ☆ Statistical modeling
- ☆ Experiment design
- ☆ Bayesian inference
- ☆ Supervised learning: decision trees, random forests, logistic regression
- ☆ Unsupervised learning: clustering, dimensionality reduction
- ☆ Optimization: gradient descent and variants

PROGRAMMING & DATABASE

- ☆ Computer science fundamentals
- ☆ Scripting language e.g. Python
- ☆ Statistical computing packages, e.g., R
- ☆ Databases: SQL and NoSQL
- ☆ Relational algebra
- ☆ Parallel databases and parallel query processing
- ☆ MapReduce concepts
- ☆ Hadoop and Hive/Pig
- ☆ Custom reducers
- ☆ Experience with xaaS like AWS



DOMAIN KNOWLEDGE & SOFT SKILLS

- ☆ Passionate about the business
- ☆ Curious about data
- ☆ Influence without authority
- ☆ Hacker mindset
- ☆ Problem solver
- ☆ Strategic, proactive, creative, innovative and collaborative


COMMUNICATION & VISUALIZATION

- ☆ Able to engage with senior management
- ☆ Story telling skills
- ☆ Translate data-driven insights into decisions and actions
- ☆ Visual art design
- ☆ R packages like ggplot or lattice
- ☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

Discover the world at Leiden University

Data Science

- Data Science =
Statistics + Computer Science + ...
- More = Better ?
 - More data = more noise = more patterns
 - Methods have to be adapted
 - Data chain



Discover the world at Leiden University

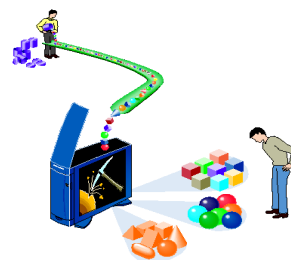
Data Science

- Data geschikt maken
 - Managen, labelen, combineren,
- Spelen met de data
 - Visualiseren, patroonherkenning, clusteren,
- Iets nuttigs mee doen
 - Classificeren, voorspellen,

Discover the world at Leiden University

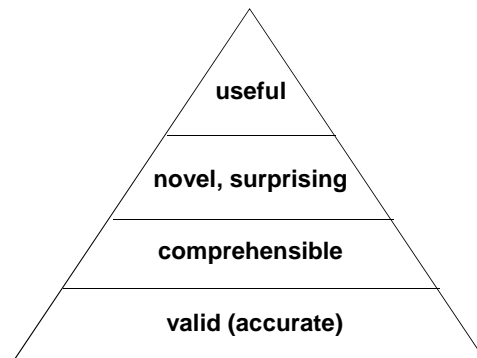
Data Mining definitions

- Secondary analysis of data
- Induction of understandable useful models and patterns from databases
- Algorithms for large quantities of data



Discover the world at Leiden University

- Data Mining is the non-trivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data



Discover the world at Leiden University

Data Mining



- Data Mining is somewhat comparable to statistics (and often based on the latter), but takes it further in the sense that whereas
 - statistics aims more at validating given hypotheses,
 - in data mining often millions of potential patterns are generated and tested, in the hope of finding some that are potentially useful.

Discover the world at Leiden University

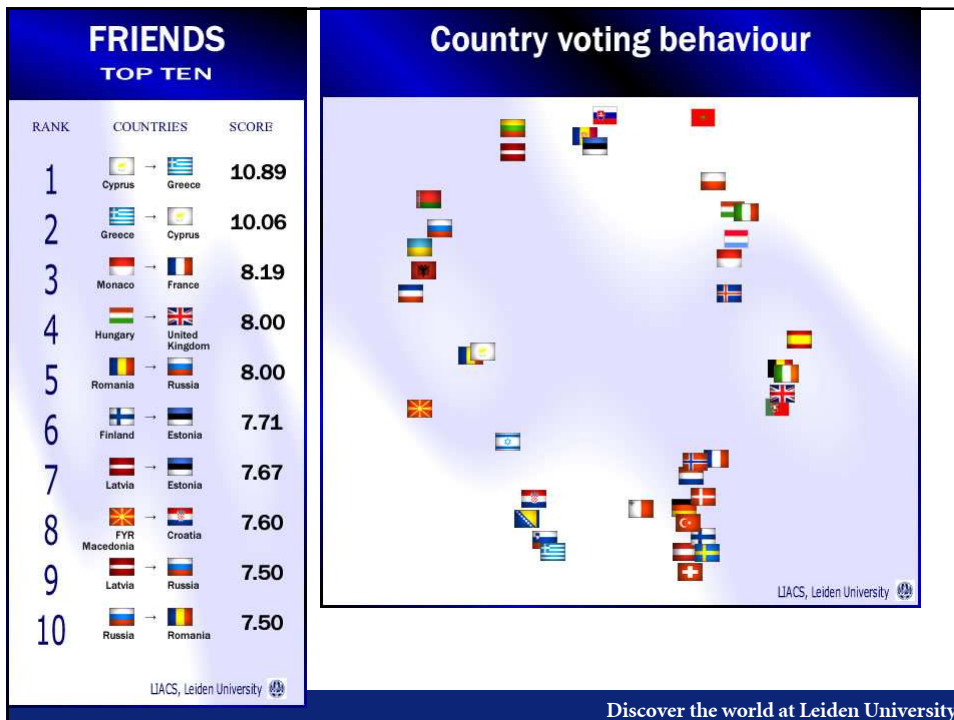
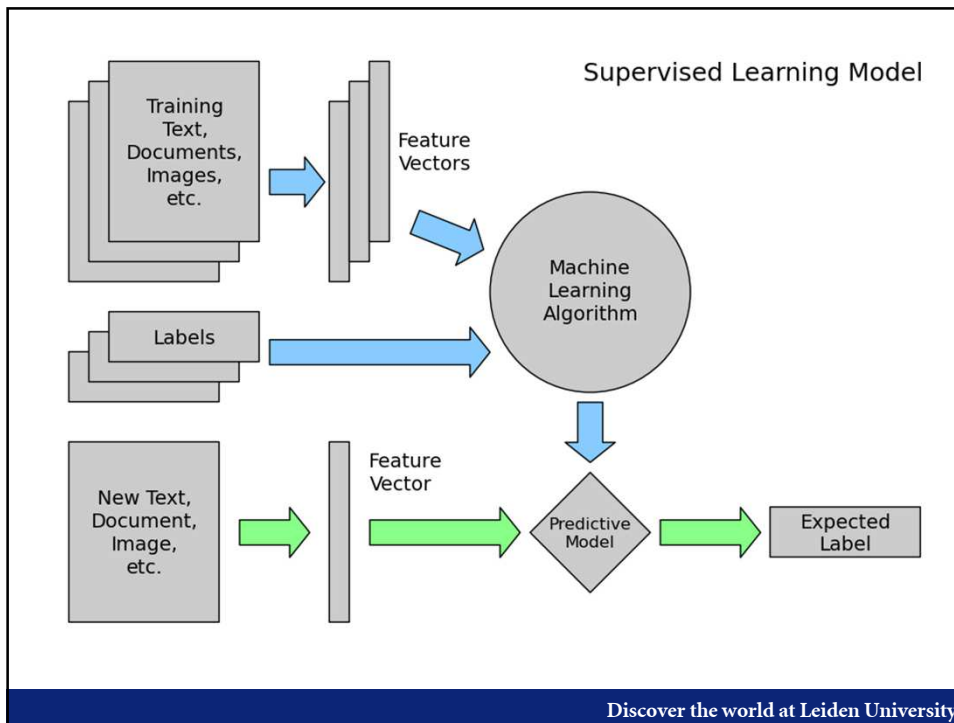
Meaningfulness of Patterns

- A risk with Big-Data mining is that an analyst can “discover” patterns that are meaningless
- Bonferroni’s principle

Meaningfulness of Patterns

Example:

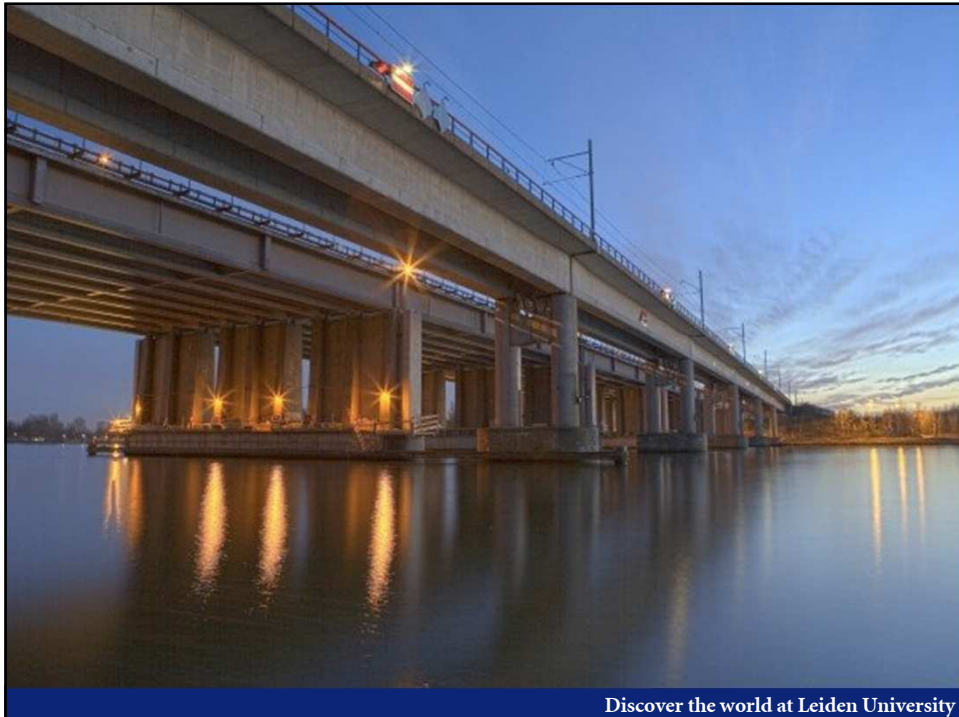
- We want to find (unrelated) people who at least twice have stayed at the same hotel on the same day
 - 10^9 people being tracked.
 - 1000 days.
 - Each person stays in a hotel 1% of the time (1 day out of 100)
 - Hotels hold 100 people (so 10^5 hotels).
 - If everyone behaves randomly (i.e., no terrorists) will the data mining detect anything suspicious?
- Expected number of “suspicious” pairs of people:
 - 250,000
 - ... too many combinations to check – we need to have some additional evidence to find “suspicious” pairs of people in some more efficient way



Monitoring

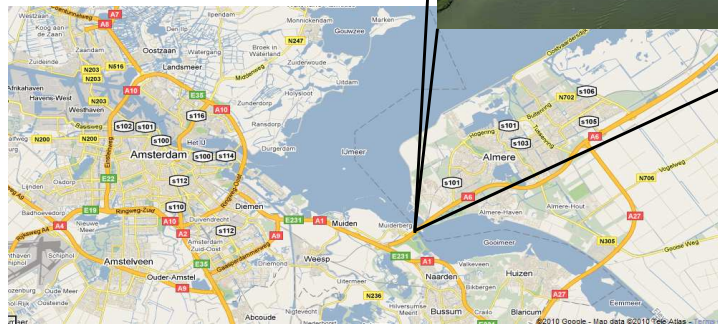
- Monitoring a Highway Bridge (Hollandse Brug)
- Health Care Data
- Sewers (SewerSense)
- School Children (Free Play)
- Animals (Oostvaardersplassen)
- Cohort Studies (Leiden 85+)
- Elderly people (Park Vossenbergh)

Discover the world at Leiden University



Discover the world at Leiden University

InfraWatch: Hollandse Brug A6 between Amsterdam and Almere

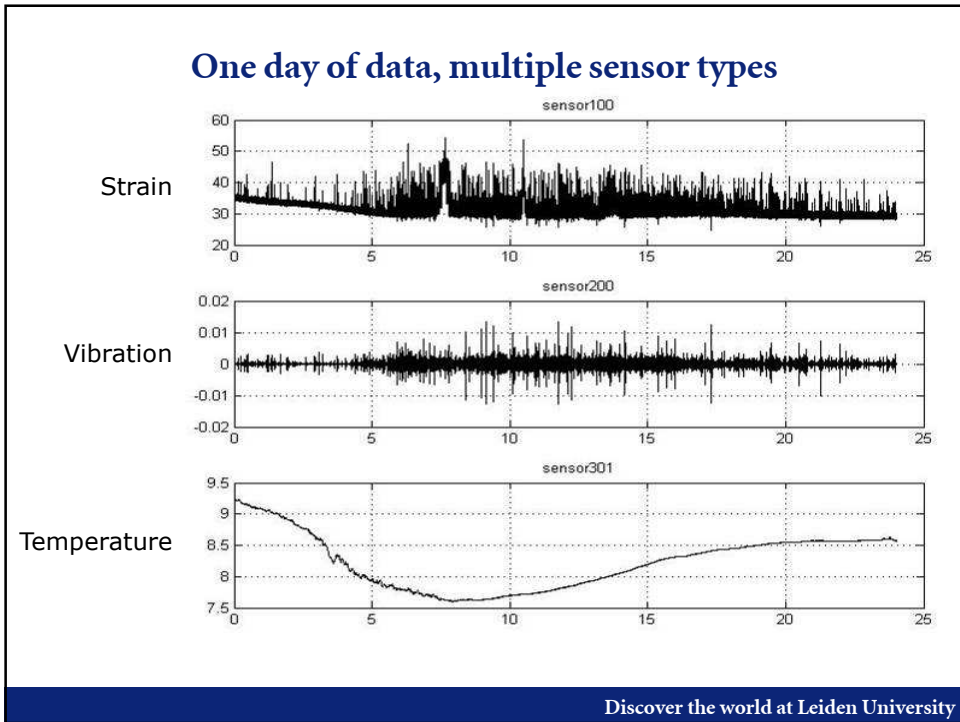
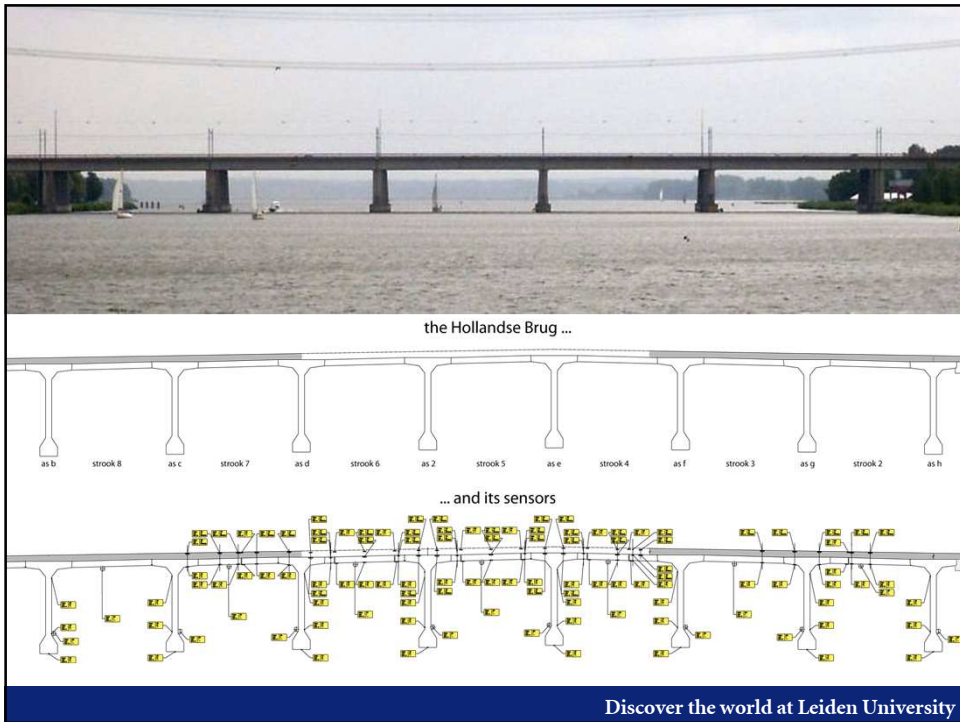


Discover the world at Leiden University

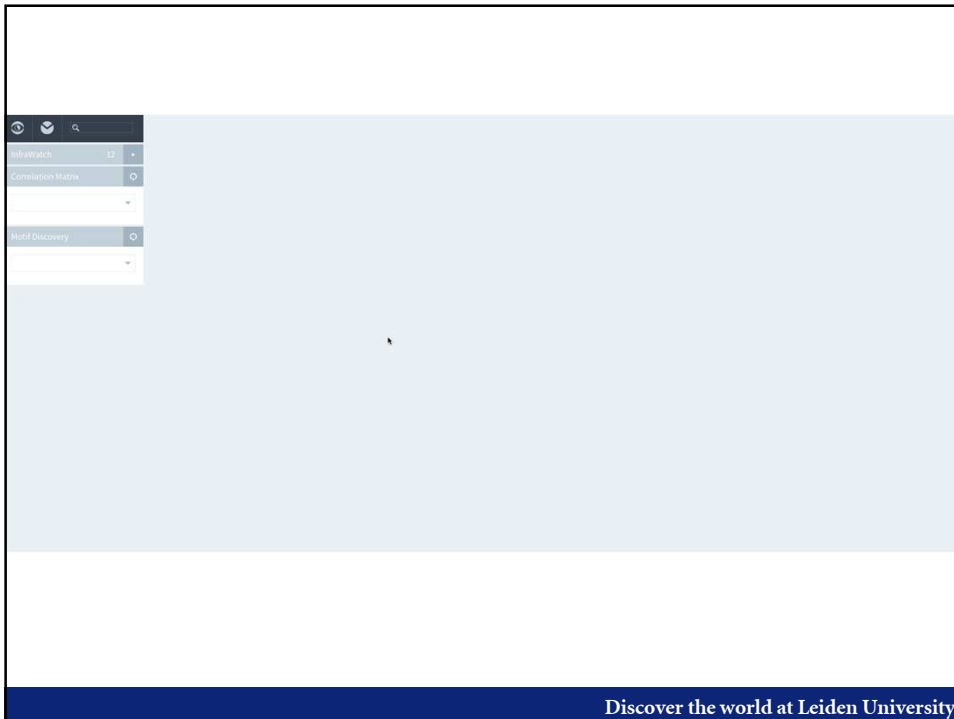
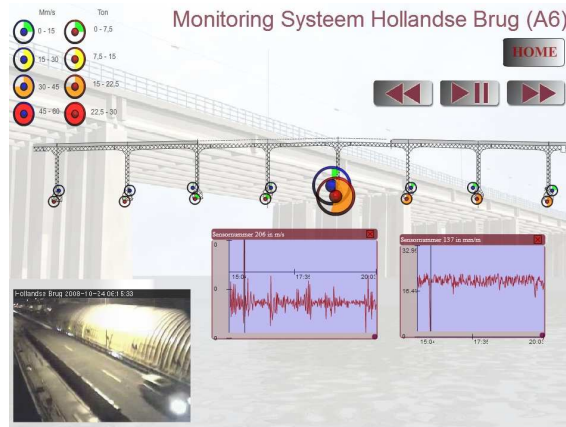
InfraWatch



Discover the world at Leiden University



Sensor Viewer



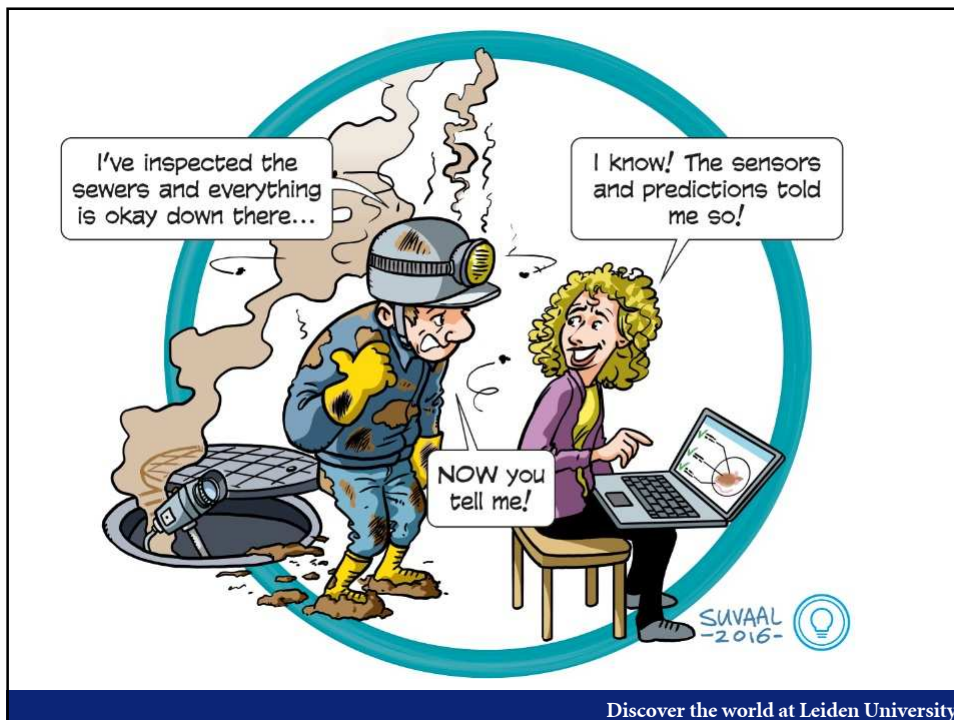


Using sensors to measure playground dynamics

18 January 2016

How do playground interactions contribute to children's social competence? Developmental psychologists Carolien Rieffe (Leiden University) and Guida Veiga (University of Évora, Portugal) joined forces with the Leiden Institute of Advanced Computer Science to investigate this. A paper on their study is currently in press.

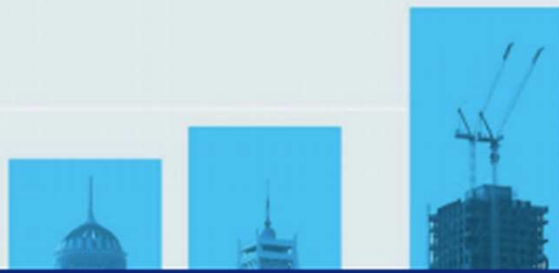
Discover the world at Leiden University



Discover the world at Leiden University

GLOBAL FRAUD REPORT

Vulnerabilities on the Rise



Discover the world at Leiden University

Data: about 2.000.000.000 records

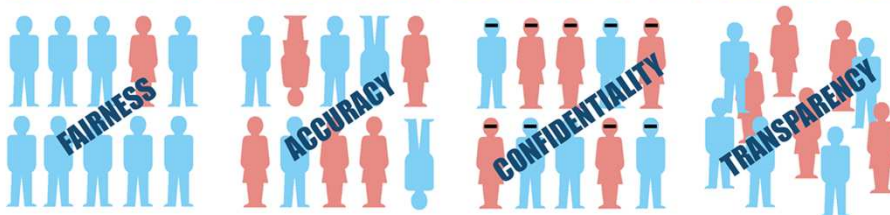
Year	Family Doctors	Dentists	Pharmacy	Mental Health Care	Fysiotherapy	Hospitals
2009				1.165.353		
2010	262.584.340	69.297.896	191.744.461	1.218.992	55.575.780	16.412.981
2011	304.654.670	68.449.999	208.515.505	1.251.854	57.068.264	17.150.880
2012	313.926.643	51.934.447	219.200.187		53.549.109	15.407.850
Totaal	881.165.653	189.682.342	619.460.153	3.636.199	166.193.153	48.971.711

Discover the world at Leiden University

	Total Amount Meuro	Hard Rules Meuro	Soft Rules Meuro	Hard and Soft Rules Meuro	Percentage %
GP	2619	15,4	6,2	21,6	0,8
Dental Care	2180	0,7	1,0	1,7	0,1
Pharmacy	5280	10,5	0,9	11,4	0,2
Mental Health	3980	4,2	-	4,2	0,2
Physiotherapy	1446	0,6	11,1	11,7	0,8
Hospitals	16676	11,9	54,7	66,6	0,4
Total	32181	43,3	73,9	117,2	0,4

Discover the world at Leiden University

RESPONSIBLE DATA SCIENCE



Discover the world at Leiden University

Responsible Data Science

- [Fairness] Data Science without prejudice: How to avoid unfair conclusions even if they are true?
- [Accuracy] Data Science without guesswork: How to answer questions with a guaranteed level of accuracy?
- [Confidentiality] Data Science that ensures confidentiality: How to answer questions without revealing secrets?
- [Transparency] Data Science that provides transparency: How to clarify answers such that they become indisputable?

Discover the world at Leiden University

Ranking of differentially expressed genes

Gene	Score
$\text{gene}_{\sigma(1)}$	score 1
$\text{gene}_{\sigma(2)}$	score 2
$\text{gene}_{\sigma(3)}$	score 3
$\text{gene}_{\sigma(4)}$	score 4
.....
$\text{gene}_{\sigma(100)}$	score 100
$\text{gene}_{\sigma(101)}$	score 101
.....
$\text{gene}_{\sigma(9905)}$	score 9905

The genes are ordered in a ranked list (e.g. according to their differential expression between the classes).

The challenge is to extract meaning from this list, to describe subgroups with common properties.

Subgroup Discovery

- Population of individuals
- A property we are interested in
- Discover the subgroups of the population that are statistically “most interesting”, i.e. are as large as possible and have the most unusual statistical (distributional) characteristics with respect to the property of interest.

Subgroup Discovery

Discovery of gene **subgroups** which

- are “higher” in the ranked list
- can be compactly summarized using
 - knowledge (GO, ENTREZ, KEGG)
 - Interactions between genes
 - ...

```

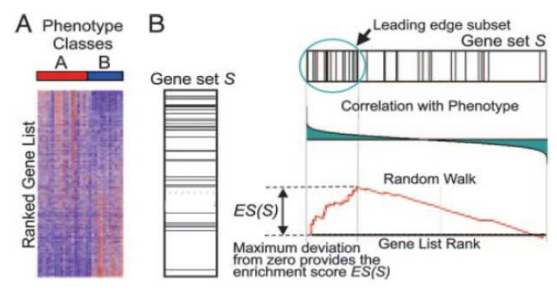
=====
Rule 1
Distribution: [10, 0]
Score: 0.807160888001244
A participants: Adh1, Coq7, Cyp2c29, Cyp2c37, Cyp3a13, Cyp3a16, Hmox1, Nudt8, Ftmt, Dpyd]
All genes in the subgroup
    have the following properties:
        cellular_component(cytoplasmic part),
        molecular_function(transition metal ion binding),
        KEGG_pathway(Metabolism of Cofactors and Vitamins),
=====

Rule 2
Distribution: [10, 0]
Score: 0.797726796885391
A participants: Atp11, Cdk5, Cdkn1b, Fxn, Glrb, Itpr3, Mecp2, Myo7a, Sod1, Sod2]
All genes in the subgroup
    have the following properties:
        biological_process(sensory perception),
        cellular_component(intracellular membrane-bounded organelle),
=====

Rule 3
Distribution: [11, 0]
Score: 0.771767724124038
A participants: Capg, Capza3, Cdkn1b, Mecp2, Mapt, Pax5, Prox1, Scin, Tmod3, Noc2l, Trim54]
All genes in the subgroup
    have the following properties:
        biological_process(regulation of organelle organization),
        cellular_component(intracellular organelle),
=====

```

Enrichment Score





Data rond Sport en Bewegen Sport Data Valley

Discover the world at Leiden University

Data Science and Sports

- Apply data science in new application domains yielding new results and insights in those domains.
- At the same time, the application domain serves as a source of inspiration for new data science research.



Discover the world at Leiden University

Vraagstukken

- Nederland staat niet in de top-10
- Kinderen zijn minder gezond
- Sociale cohesie neemt af
- Fraude en matchfixing
- Mensen bewegen te weinig
- Genezen door bewegen

Discover the world at Leiden University



Discover the world at Leiden University

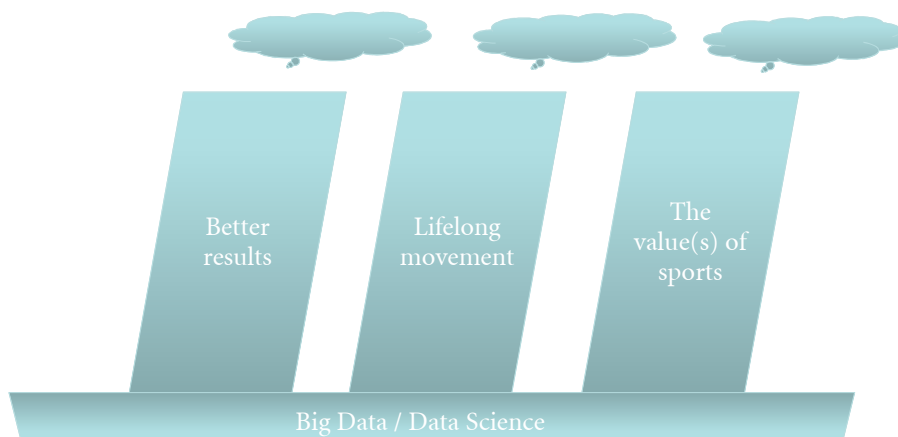
Het kabinet verdubbelt de in 2016 afgesproken structurele intensivering voor de topsport van 10 miljoen euro per jaar naar 20 miljoen euro per jaar om meer kansen te bieden aan onze Olympische en Paralympische teams. Daarnaast komt er meer ruimte voor topsport talenten om onderwijs en topsport combineren.

Het kabinet trekt structureel 5 miljoen euro extra uit ter ondersteuning van de organisatie van sportevenementen in Nederland, waaronder EK's, WK's en multisportevenementen. Het initiatief voor de organisatie van een evenement ligt altijd bij de sport en haar partners.

Het kabinet zal het gesprek met de bonden aangaan over de handhaving van de openbare orde bij evenementen en de omgang met gedragingen tijdens evenementen en risicowedstrijden.

Het kabinet gaat steviger inzetten op de aanpak van dopinggebruik, matchfixing, corruptie en misbruik in de sport.

Discover the world at Leiden University



**Knowledge Agenda
Sports and Movement**

Discover the world at Leiden University

Advies
op Maat

Tailored
Advice



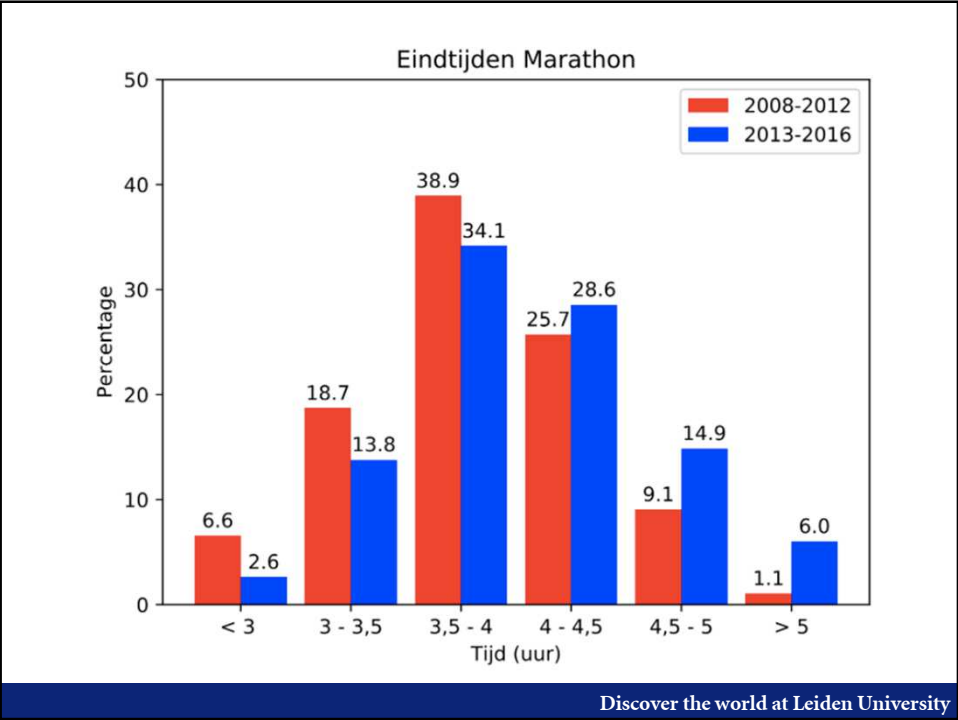
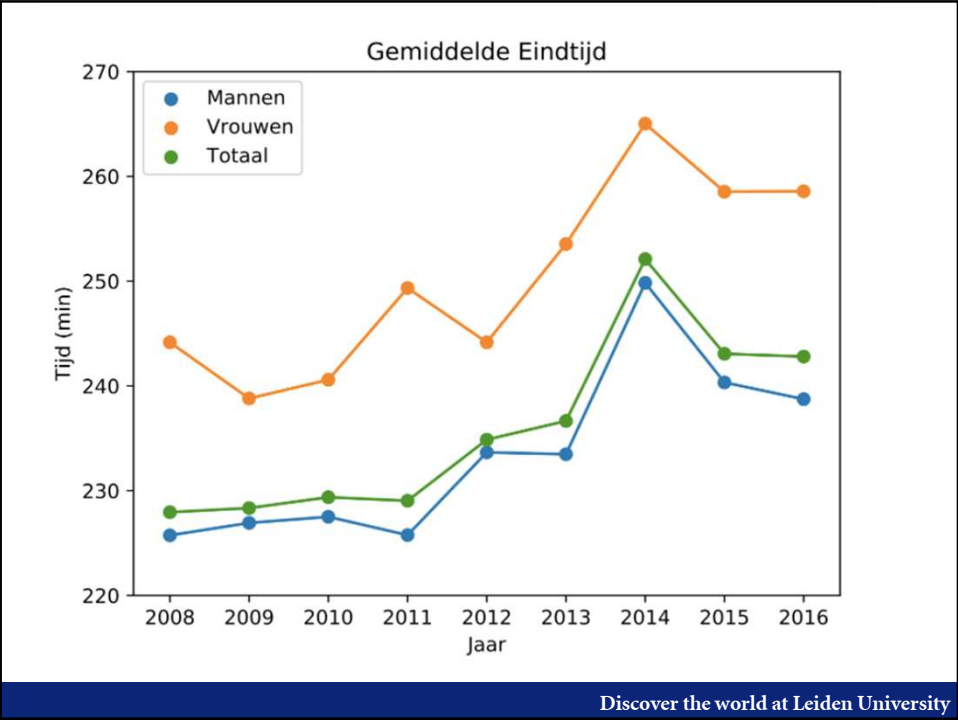
Discover the world at Leiden University

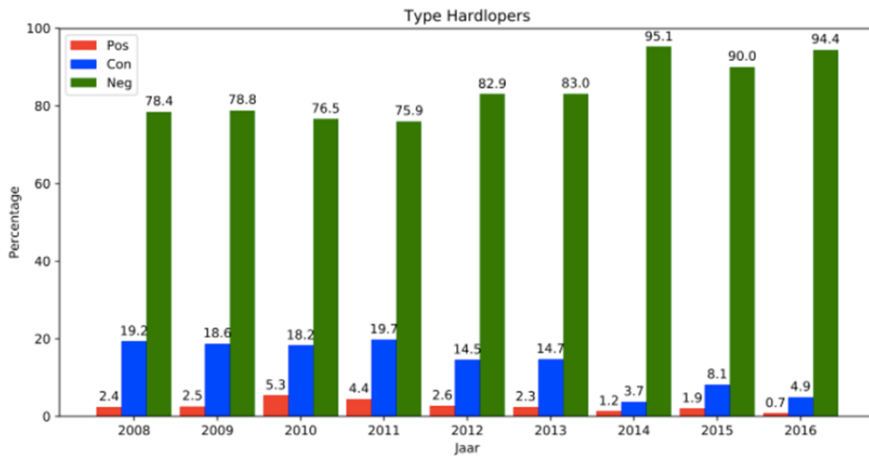
Healthy
Movement

Joost
Kok



Discover the world at Leiden University



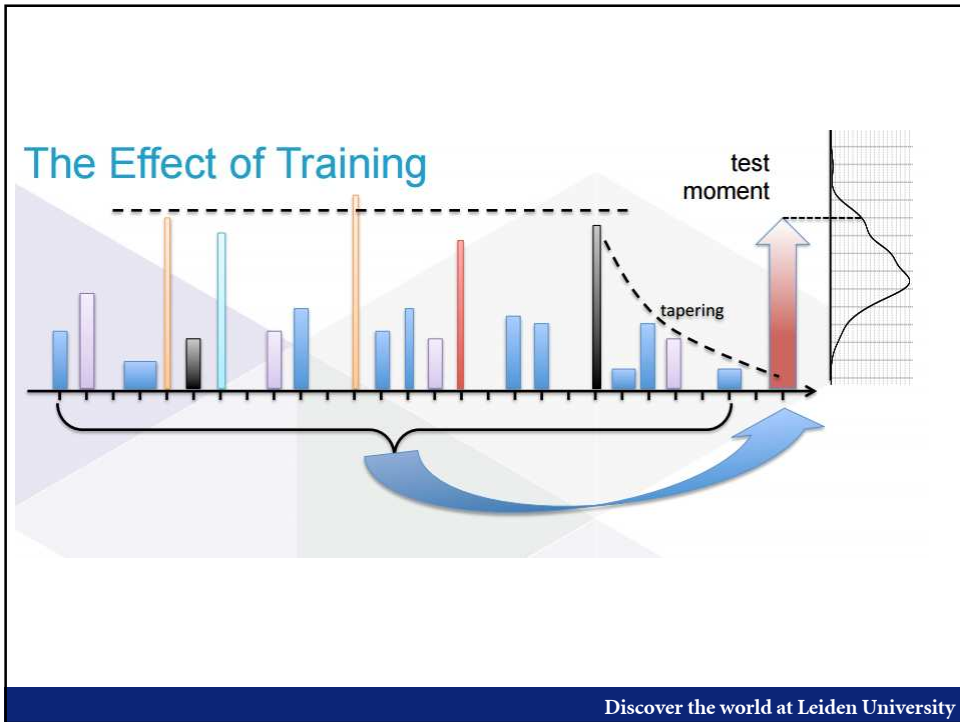


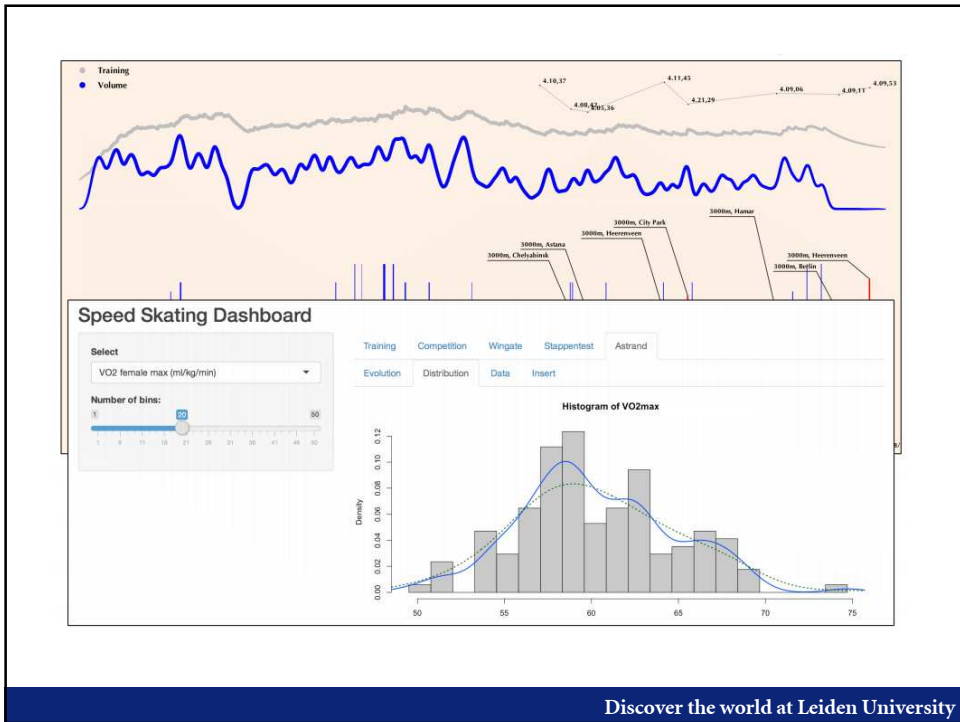
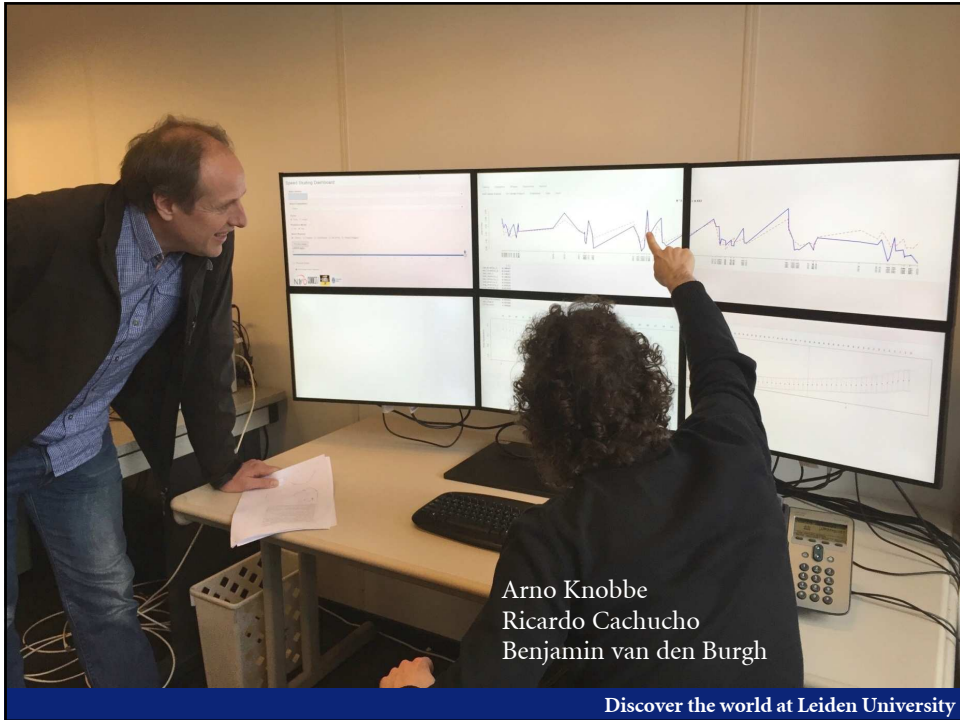
pos. : lopers die de tweede helft sneller lopen dan eerste helft.
 con. : lopers die een constant tempo lopen.
 neg. : lopers die de tweede helft langzamer lopen dan eerste helft.

Loop Data

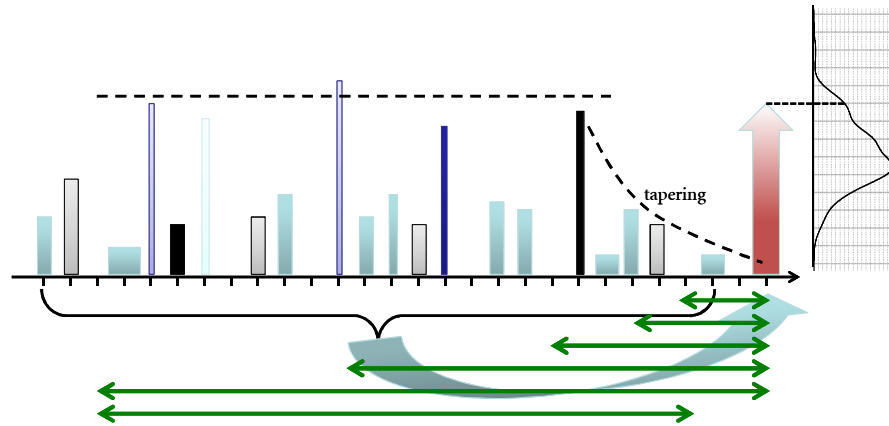
- Blessurepreventie
- Strategie







The Effect of Training



Discover the world at Leiden University

Training parameters

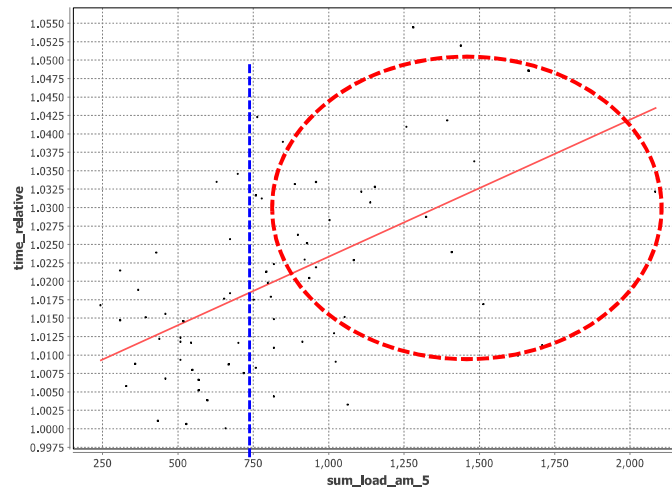
- Hundreds of training parameters that capture various aspects of periodization

- Sum of duration over 14-day period
- Max of intensity over 2-day period
- Sum of duration over 21-day period, morning sessions
- Sum of duration over 21-day period, intensities 6, 7, 8, 9
- Maximum of load over 7-day period, cycling
- Count of cycling sessions over 4-day period

etcetera

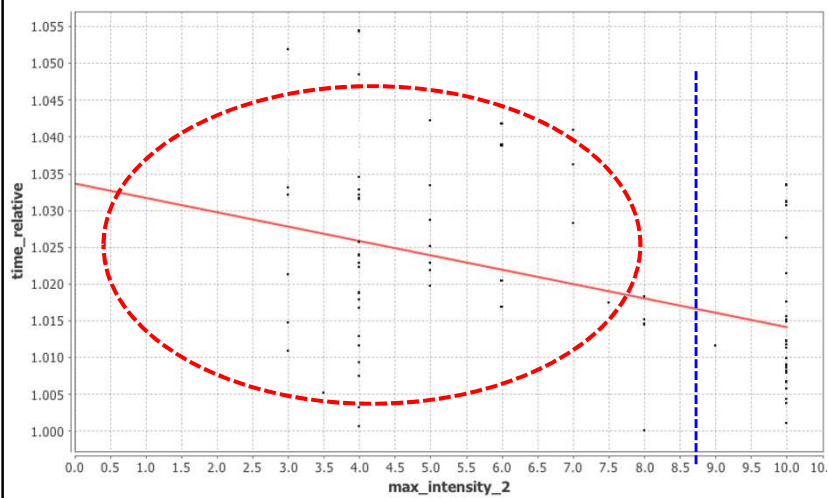
Discover the world at Leiden University

Sum of loads over last 5 days, morning sessions



Discover the world at Leiden University

Maximum intensity 8 or higher in the two days prior to the race



Discover the world at Leiden University

Some Findings

- To increase *aerobic capacity*, make sure you
 - include at least one exercise longer than 3.5 hours
 - ...over the period of 14 to 3 days before the test moment
 - avoid loads above 240 in the mornings, 2 days window

VO₂max will increase by 4%

- total time in intensity zone [1, 4] above 850 min/w, 21 days window
- average intensity above 3.8, 14 days window

VO₂max will increase by 11%



Discover the world at Leiden University

Results

- Longitudinal, detailed data, great potential
- Actionable results, taken up by coach
- No major revolution in training:
 - subtle tweaking of parameters
- Joint scientific publication
- Team funding

Sports Analytics for Professional Speed Skating

Arno Knobbe · Jac Orié · Nico Hofman · Benjamin van der Burgh · Ricardo Cachucho

the date of receipt and acceptance should be inserted later

Abstract Elite athletes need to optimise every detail of their training routines, if they want to compete at the top of their sport. Training schedules are becoming increasingly complex, and a large number of parameters of these schedules need to be tuned to the specific physique of a given athlete. In this paper, we describe how extensive analysis of historical data can help optimise these para-

Investment Decisions

papendal
hotel • sport • congres

Marvin Meeng

Discover the world at Leiden University

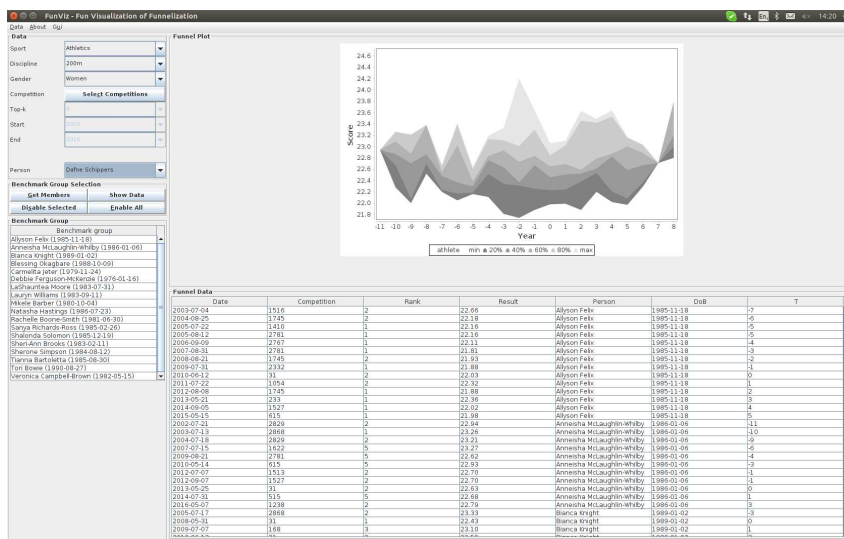
Discover the world at Leiden University



Discover the world at Leiden University

Funnels

Marvin Meeng



Discover the world at Leiden University

FunViz - Fun Visualization of Funnellization

Data about Gg

Sport: Athletics
 Discipline: 100m
 Gender: Men
 Competition: Select Competitions
 Top-K: 10
 Start: 2000
 End: 2016
 Person: Churandy Martina

Benchmark Group Selection
 Get Members Show Data
 Disable Selected Enable All

Benchmark Group

Benchmark group
 Andre De Grasse (1994-11-10)
 Asafa Powell (1982-11-23)
 Ato Boldon (1973-12-30)
 Bernard Williams (1978-01-19)
 Darrel Brown (1984-10-11)
 Darren Campbell (1973-09-12)
 Demick Atkins (1984-01-05)
 Francis Obikwelu (1978-11-22)
 Justin Gatlin (1982-02-10)
 Kim Collins (1976-04-05)
 Maurice Greene (1974-07-23)
 Michael Frater (1982-10-06)
 Nesta Carter (1983-10-11)
 Obadiah Thompson (1976-03-30)
 Richard Thompson (1965-06-17)
 Tyson Bromell (1985-07-10)
 Tyson Gay (1982-08-09)
 Usain Bolt (1986-08-21)

Date	Competition	Rank	Result	Person	DoB	T
2015-08-23	2781	3	9.92	Andre De Grasse	1994-11-10	1
2014-05-28	1822	8	10.05	Andre De Grasse	1994-11-10	1
2003-09-05	1527	1	10.02	Asafa Powell	1982-11-23	4
2004-09-03	1527	1	9.97	Asafa Powell	1982-11-23	13
2005-06-14	167	1	9.77	Asafa Powell	1982-11-23	12
2006-06-11	259	2	9.73	Asafa Powell	1982-11-23	1
2006-08-18	2717	1	9.77	Asafa Powell	1982-11-23	13
2007-09-09	2019	1	9.78	Asafa Powell	1982-11-23	0
2008-09-02	168	1	9.72	Asafa Powell	1982-11-23	1
2009-08-16	2781	3	9.84	Asafa Powell	1982-11-23	2
2010-06-04	263	1	9.72	Asafa Powell	1982-11-23	3
2011-06-30	168	1	9.78	Asafa Powell	1982-11-23	4
2012-06-07	263	2	9.85	Asafa Powell	1982-11-23	5
2013-07-04	168	1	9.88	Asafa Powell	1982-11-23	6
2014-09-07	2019	2	9.90	Asafa Powell	1982-11-23	7
2015-07-04	1518	1	9.83	Asafa Powell	1982-11-23	8
2016-05-28	1923	2	9.94	Asafa Powell	1982-11-23	9
1992-09-01	1745	44	10.17	Ato Boldon	1973-12-30	4
1995-08-08	2781	3	10.03	Ato Boldon	1973-12-30	5
1996-07-27	1745	4	9.86	Ato Boldon	1973-12-30	4
1997-08-03	2781	5	10.02	Ato Boldon	1973-12-30	3
1998-09-17	515	1	9.88	Ato Boldon	1973-12-30	2
2000-09-23	1745	2	9.99	Ato Boldon	1973-12-30	0
2001-07-13	263	1	9.88	Ato Boldon	1973-12-30	1
2004-08-22	1745	45	10.01	Ato Boldon	1973-12-30	4
1999-07-25	1822	1	10.08	Bernard Williams	1978-01-19	12
2000-08-25	1527	1	10.01	Bernard Williams	1978-01-19	11
2001-08-05	2781	2	9.94	Bernard Williams	1978-01-19	0
2002-08-18	1822	2	10.08	Bernard Williams	1978-01-19	1

Discover the world at Leiden University

Statistical Methods for Ranking Data

Springer

Discover the world at Leiden University



Discover the world at Leiden University

SPORT DATA CENTER MAP ATHLETES COUNTRIES PREDICTIONS ABOUT

Choose discipline(s): Horizontal Bar Men Rating History Rating by Date Age potential winners athletes overview info About

Choose Athlete(s): Epke Zonderland

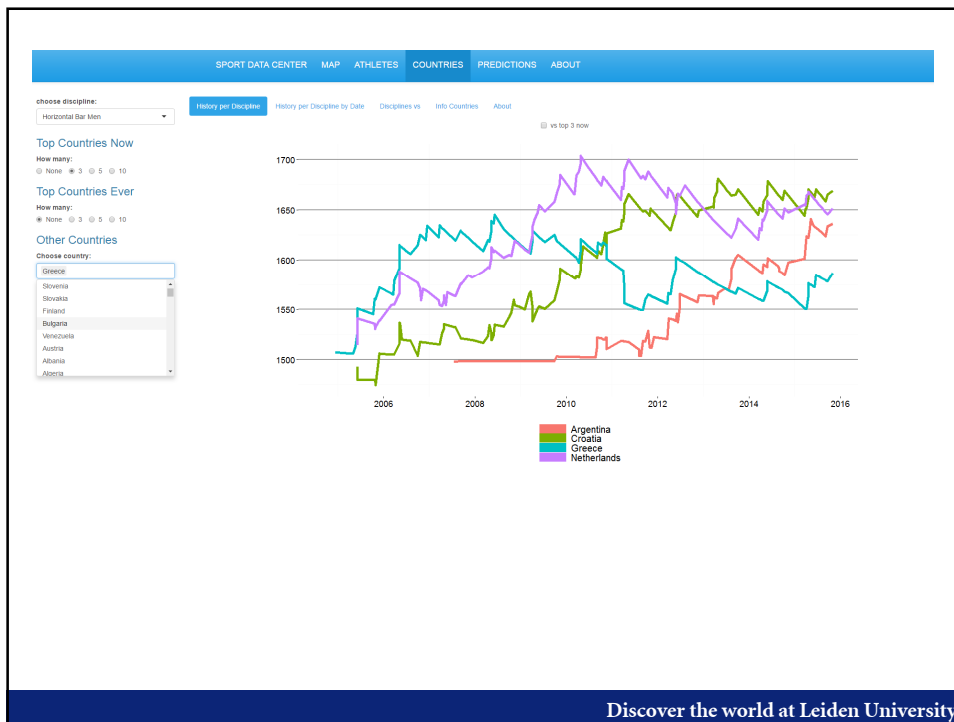
Epke Zonderland

Enve Arabacoglu
Endrasi
Enkhmurkh Munkhjargal
Enrico Pozzo
Enzo Navaro
Enzo Bernardoni
Epke Zonderland

Personal information for Epke Zonderland

Name: Epke Zonderland
Discipline: Horizontal Bar
Category: Men
Age: 29.5
Country: Netherlands
ELO: 1651.31
Rank: 2

Discover the world at Leiden University



Challenges

- Develop Elo-like ratings
 - Probabilities
 - Head-to-head including differences
 - Doubles
 - Teams

Discover the world at Leiden University



Discover the world at Leiden University

- Analyse opponent tactics
- Detect strengths/weaknesses in strategy
- Automatic game plans
- Serious games / training
- Player scouting
- Improved media coverage

Discover the world at Leiden University

Research question

- Using an algorithmic approach that analyses possession switches during football matches, which variables have the highest influence on these moments?

Discover the world at Leiden University

Discover the world at Leiden University

Data

Available data

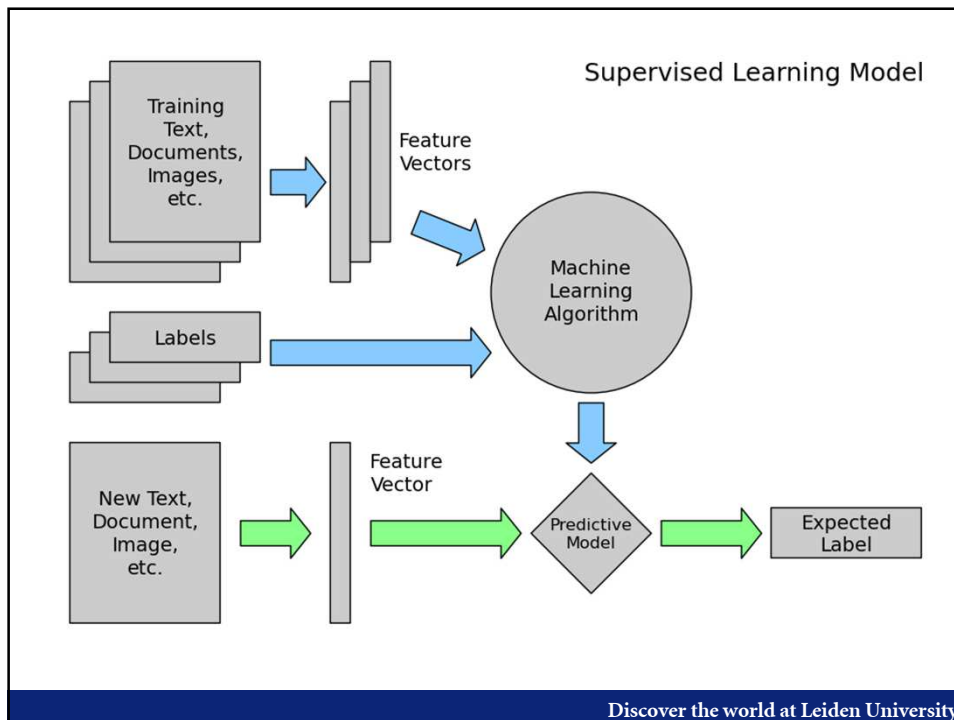
- Positional data
- Event data

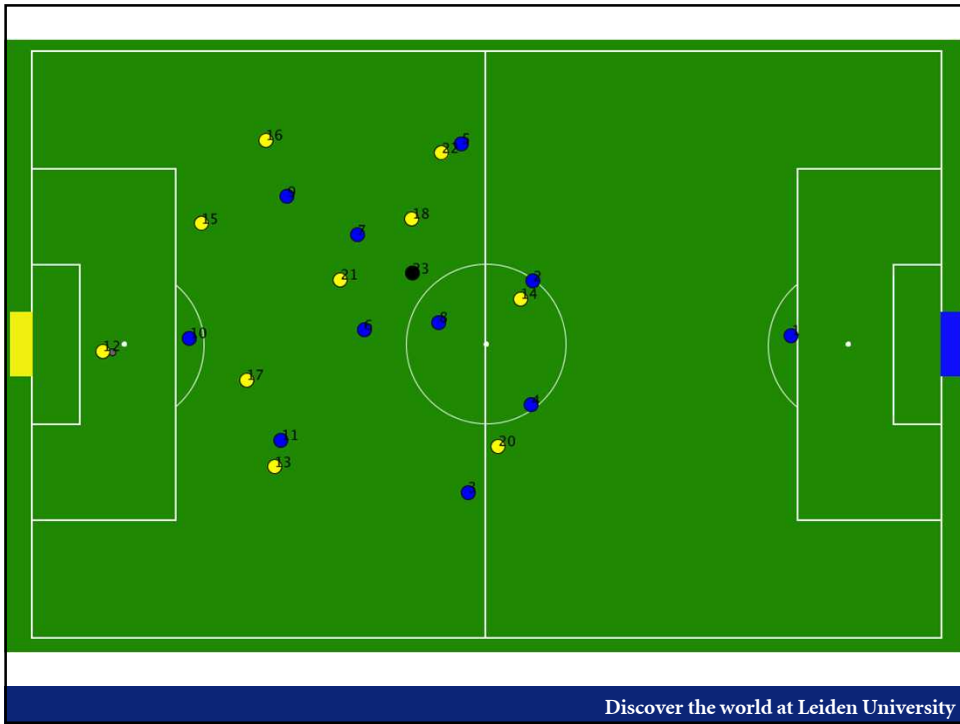
Timestamp (ms)	X	Y	Name
40,000	17.027	5.364	Home Team Player
40,000	-4.051	-3.378	Visiting Team Player 1
40,000	99.999	99.999	Visiting Team Player 2
40,000	11.094	-10.509	Referee
40,000	38.09	-7.985	Ball

Table 3.1: Frame 40,000 (first frame of 40th second).

Time(ms)	Half	Category	Player	Team	Attribute	Definition	Match
7,132	1	attacking action	Player 8	Home Team	head — duel touched	isPossessionLoss — isDuelPart — isDuelAir — isAerial	Match 1
10,905	1	pass	Player 15	Visiting Team	right foot	isPassCompleted — isPassWide	Match 1
11,930	1	pass	Player 16	Visiting Team	right foot — direct	isPossessionLoss — isPassForward — isPassLong	Match 1
14,931	1	pass	Player 3	Home Team	direct — left foot	isPossessionGain — isPossessionGainInterception — isPassCompleted — isPassForward	Match 1
17,030	1	attacking action	Player 4	Home Team	duel touched	isDuelPart — isDuelWonByAttacker — isDuelStanding	Match 1

Table 3.2: Event data example.





Sum of distances to ball

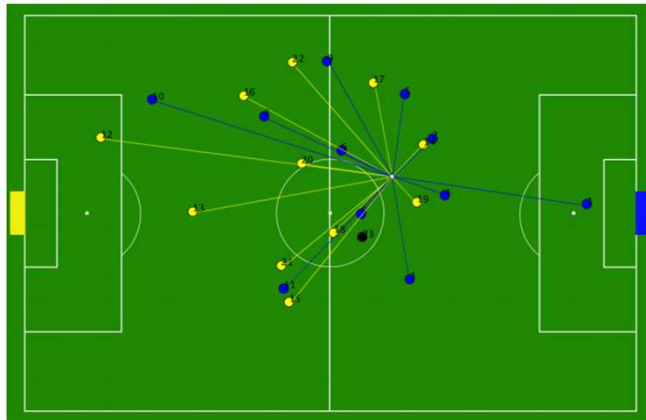


Figure 4.1: Distance to the ball.

Sum of Distances to closest opponent

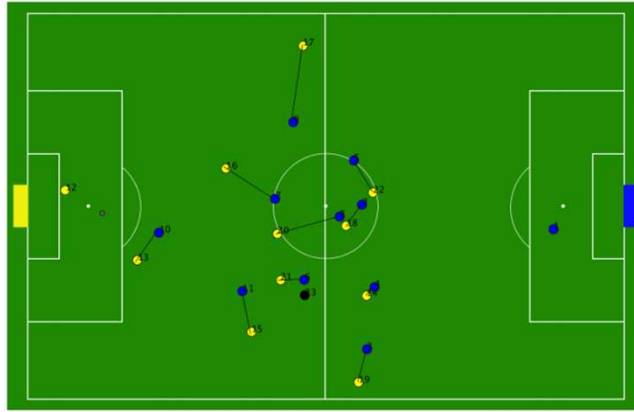
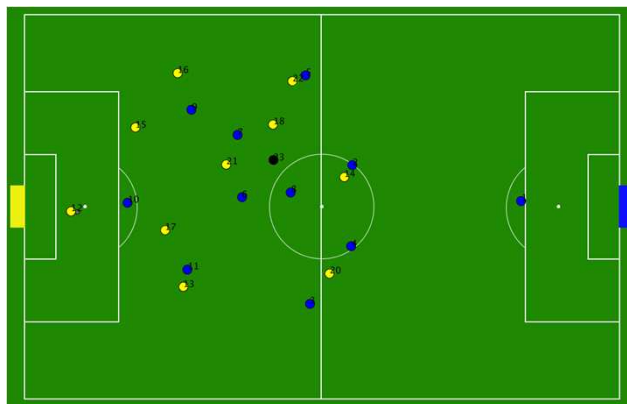


Figure 4.2: Optimal direct opponent.

Discover the world at Leiden University

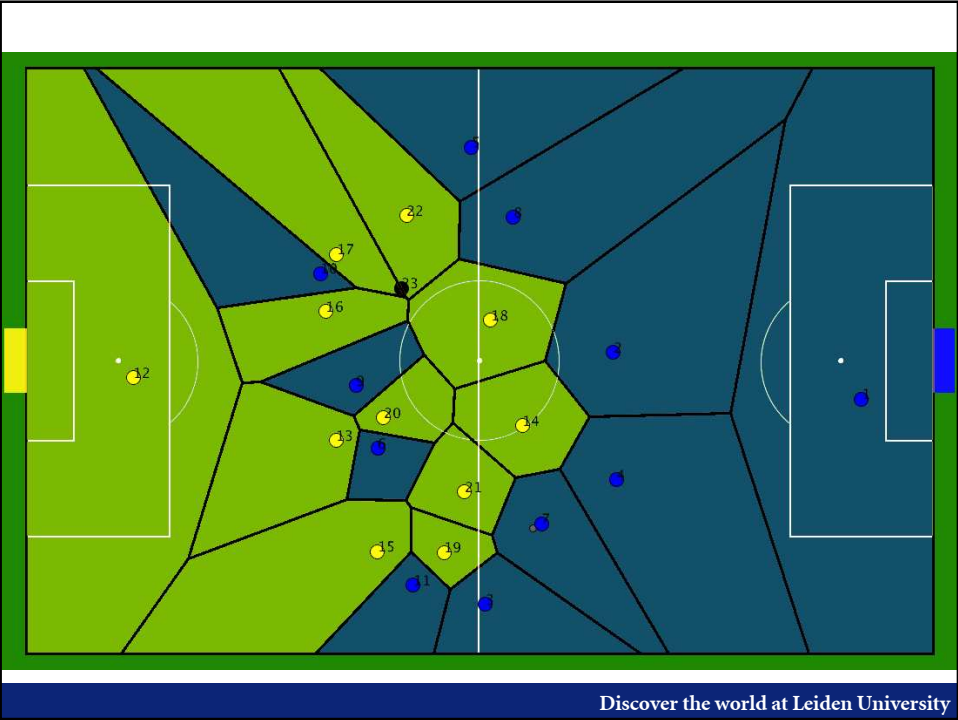
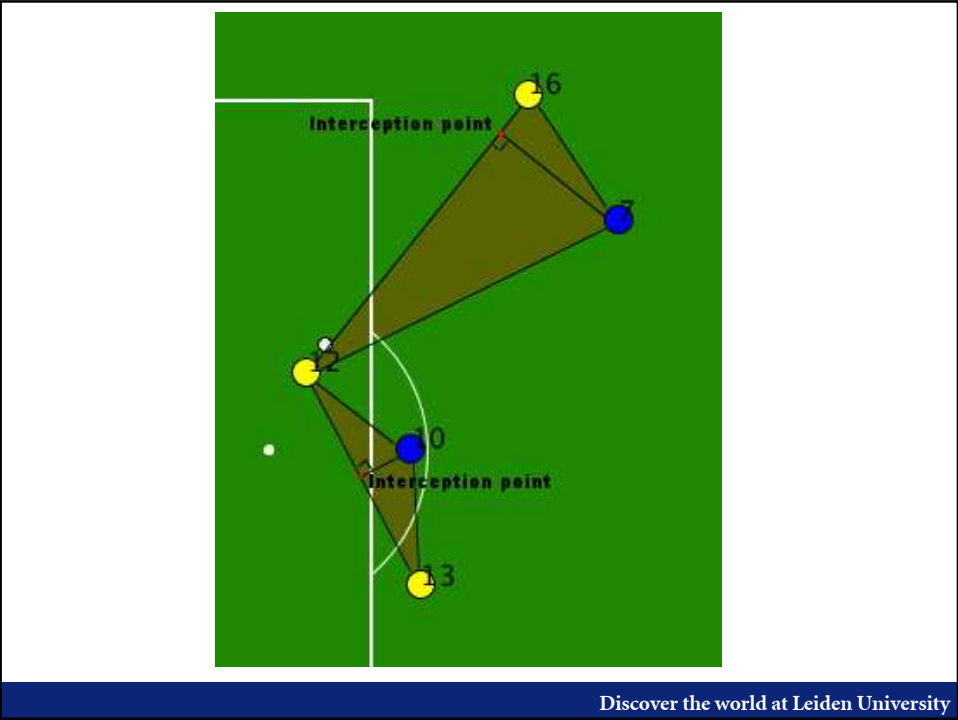
Counts

- # Players of own team within X meters
- # Players of opposite team within X meters



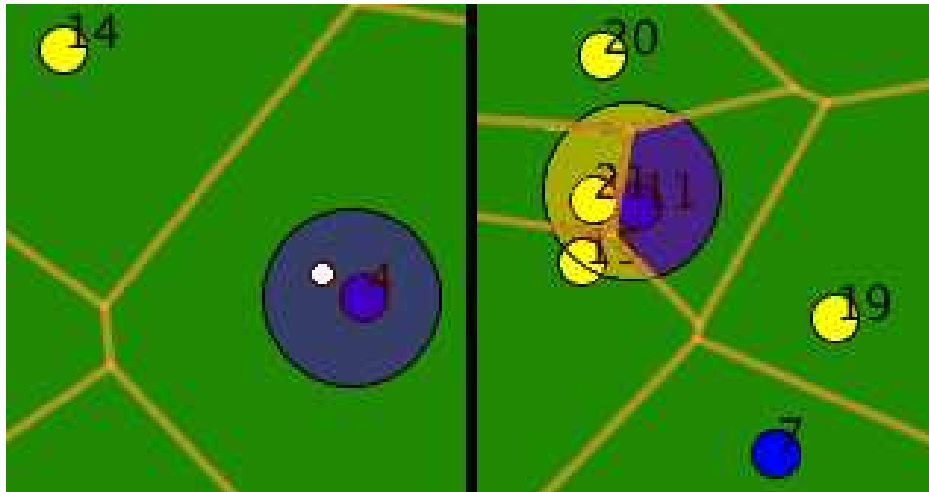
Discover the world at Leiden University

Discover the world at Leiden University



Voronoi Diagrams

- <https://www.youtube.com/watch?v=7eCrHAv6sYY>
- <https://www.youtube.com/watch?v=Y5X1TvN9TpM>
- <https://www.youtube.com/watch?v=k2P9yWSMaXE>



Features

Football knowledge

- Possession switch
- Pressure

Features

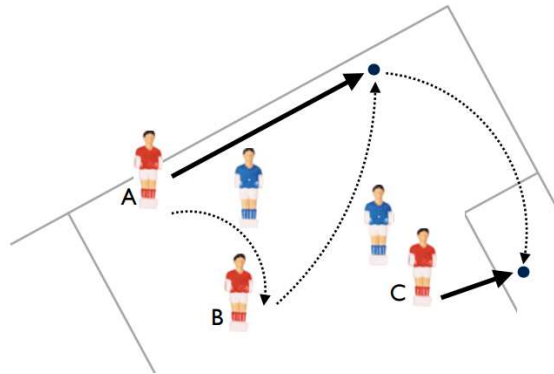
- Distances
- Positions
- Counts
- Ratios
- Surfaces
- Possible passes

Subgroup Discovery

- Pattern

The distance between the ball and the closest 3 players of the defending team should be smaller or equal to 27 meters for a possession switch to occur

- Pattern = “interesting” event
- E.g., A plays 1-2 with B and crosses to C



Slides by Ulf Brefeld

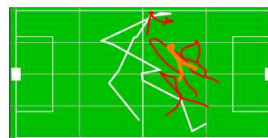
Discover the world at Leiden University

- Individual level
- Group level
- Team level

- 4 defence players
→ game initiations



- 4 offence players
→ scoring opportunities



Different Scales

Discover the world at Leiden University

Data is de nieuwe olie

Kennistransfer van

- extern naar intern
- topsport naar breedtesport
- sport naar sport
- expert naar groot publiek



Advies op maat

- van groep naar individu
- van lid naar klant

Data is the new oil.
It's only useful when
it's refined!

Jess Greenwood, Contagious

Discover the world at Leiden University

Fraud & Risks



Discover the world at Leiden University



Data rond Sport en Bewegen

Discover the world at Leiden University



- Datawetenschappers
- Sport Data Valley

- Domeinkennis essentieel
- Data & sport field labs



Universiteit Leiden



LEIDEN UNIVERSITY MEDICAL CENTER



Discover the world at Leiden University



Vijf onderzoekslijnen in sportdata analyse

1. Topsport
2. Breedtesport
3. Aangepast sporten
4. Economische waarde van sport
5. Fraude en risico's

Discover the world at Leiden University

Data Science and Sports

- Apply data science in new application domains yielding new results and insights in those domains.
- At the same time, the application domain serves as a source of inspiration for new data science research.



Discover the world at Leiden University

Sportdata

- Data Scientists have no data
- Positive, easy to explain
- Quantified Self
- Testbed



Discover the world at Leiden University



Discover the world at Leiden University