

Data Science
Spyros Samothrakis
Deputy Director, IADS
University of Essex

February 12, 2019

Data Science

Analytics

Examples

Prescriptive Analytics

BETTER SCIENCE THROUGH DATA

Hey, Tony, Stewart Tansley, and Kristin M. Tolle. “Jim Gray on eScience: a transformed scientific method.” (2009).

- ▶ Thousand years ago: empirical branch
 - ▶ You observed stuff and you wrote down about it
- ▶ Last few hundred years: theoretical branch
 - ▶ Equations of gravity, equations of electromagnetism
- ▶ Last few decades: computational branch
 - ▶ Modelling at the micro level, observing at the macro level
- ▶ Today: data exploration
 - ▶ Let machines create models using vast amounts of data

MIXING STATISTICS, PHILOSOPHY OF SCIENCE AND MACHINE LEARNING

- ▶ Wu, C. F. J. “Statistics= data science.” (1997).
- ▶ Breiman, Leo. “Statistical modeling: The two cultures (with comments and a rejoinder by the author).” *Statistical Science* 16.3 (2001): 199-231.
- ▶ Science is the epistemology of causation
- ▶ Data science is often science on observational data
 - ▶ But quite often we only care about predictions
- ▶ Possibly a re-branding of data mining, machine learning, artificial intelligence, statistics

BETTER BUSINESS THROUGH DATA

- ▶ There was that report by Mckinsey

Manyika, James, et al. “Big data: The next frontier for innovation, competition, and productivity.” (2011).

- ▶ Urges everyone to monetise “Big Data”
- ▶ Use the data provided within your organisation to gain insights
- ▶ Has some numbers as to how much this is worth
- ▶ Proposes a number of methods, most of them associated with machine learning and databases

MORE IS DIFFERENT

- ▶ Anderson, Philip W. “More is different.” Science 177.4047 (1972): 393-396.
- ▶ The idea of emergence
- ▶ You put stuff together, you go from physics to chemistry
 - ▶ ...from chemistry to biology
 - ▶ ...from biology to psychology and zoology
 - ▶ ...from psychology to sociology

quantity changes into quality

OVERVIEW - ANALYTICS

- ▶ Descriptive
 - ▶ “I’ll create high level, compressed views of your data”
 - ▶ What is sometimes termed analytics, unsupervised learning
- ▶ Predictive
 - ▶ “I’ll try to predict a possible version of your future”.
 - ▶ Show me some good customers in your database, I’ll try getting you more
 - ▶ Supervised learning (mostly)
- ▶ Prescriptive
 - ▶ “What should I do to achieve certain results?”
 - ▶ **Reinforcement Learning, bandits and causality**
 - ▶ e.g. Send the right sequence of e-mails to engage an audience

¹Chris Wiggins. "Lectures delivered Aug 8-9, 2016 at MLSS.cc (Arequipa, Peru).

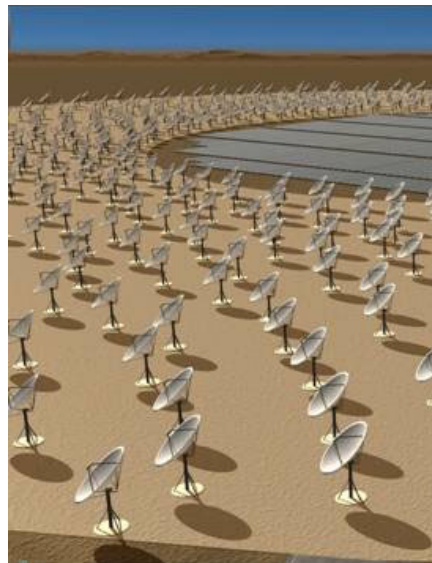
<https://www.slideshare.net/chriswiggins/machine-learning-summer-school-2016/75>

WAAAAAAIT...

- ▶ “We’ve been doing this thing for years”
 - ▶ econometrics
 - ▶ statistics
 - ▶ you-name-it
- ▶ Mostly true, in fact most of these ideas have been developed in parallel in many communities
 - ▶ Deep insights, larger scale

CLASSIC SCIENCE

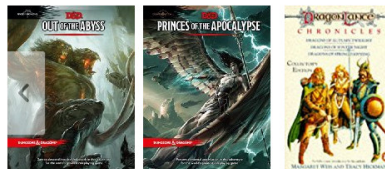
- ▶ The original data science field
- ▶ SKA (The Square Kilometer Array) ~ 4.6 EB expected (i.e. 4.6×10^6 TB), (Zhang, Yanxia, and Yongheng Zhao. “Astronomy in the Big Data Era.” Data Science Journal 14 (2015).)¹
- ▶ Bioinformatics
- ▶ Medical science



¹<http://datascience.codata.org/article/10.5334/dsj-2015-011>

RECOMMENDER SYSTEMS

- ▶ One of the most popular applications of data science
- ▶ Propose products to customers based on past history
- ▶ Almost all online vendors do it
- ▶ Made popular by the Netflix prize

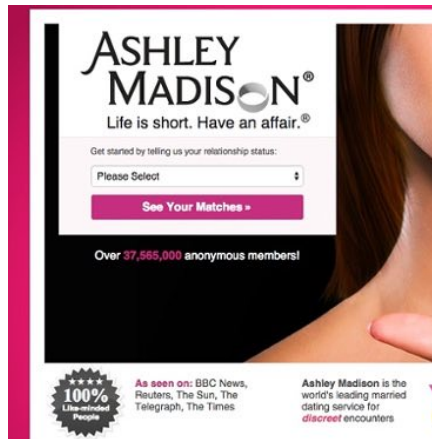


Digital Cameras best sellers [See more](#)



DATA JOURNALISM

- ▶ Wikileak style data dumps are everywhere
- ▶ The Ashley-Madison Affair, 2015
- ▶ “Just three in every 10,000 female accounts on infidelity website are real”
- ▶ “The website claims 5.5 million of its 37 million accounts are ‘female’ ”



²<http://www.independent.co.uk/life-style/gadgets-and-tech/news/ashley-madison-hack-just-three-in-every-10000-female-accounts-on-infidelity-website-are-real-10475310.html>

FINANCE & INSURANCE

- ▶ Predict stock prices
- ▶ Insurance models
- ▶ Credit scores
- ▶ Sudden drops in share prices sometimes caused by defective algorithms



³<http://www.bbc.com/news/business-34264380>

POLITICS (CURRENT)

“...This included a) integrating data from social media, online advertising, websites, apps, canvassing, direct mail, polls, online fundraising, activist feedback, and some new things we tried such as a new way to do polling (about which I will write another time) and b) having experts in physics and machine learning do proper data science in the way only they can – i.e. far beyond the normal skills applied in political campaigns...”

Dominic Cummings's (Head of *Vote Leave*) Blog²

⁴<https://dominiccummings.wordpress.com/2016/10/29/on-the-referendum-20-the-campaign-physics-and-data-science-vote-leaves-voter-intention-collection-system-vics-now-available-for-all/>

POLITICS (HISTORICAL)

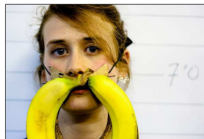
- ▶ New Yorker - THE PLANNING MACHINE: Project Cybersyn and the origins of the Big Data nation³
- ▶ Cybersyn / Chile during Alliente's rule, co-designed by Stafford Beer
- ▶ Plan was to use data fed directly from each industry to automate production



⁵<http://www.newyorker.com/magazine/2014/10/13/planning-machine>

QUESTION ANSWERING

- ▶ e.g. Antol, Stanislaw, et al. “VQA: Visual question answering.” Proceedings of the IEEE International Conference on Computer Vision. 2015.⁴
- ▶ Input can be videos, websites, et
- ▶ Think google



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



Is this person expecting company?
What is just under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?

⁶http://www.cv-foundation.org/openaccess/content_iccv_2015/papers/Antol_VQA_Visual_Question_ICCV_2015_paper.pdf

DIGITAL MARKETING

- ▶ Is a new product I just created well received by our customers?
- ▶ Is a new marketing campaign e-mail sent detrimental to our efforts?
- ▶ What is the content a chain of e-mails should have?
- ▶ What adverts should I present to a user?

BUSINESS ANALYTICS

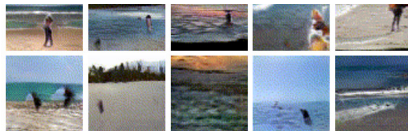
- ▶ Churn models
 - ▶ Why are my customers leaving?
- ▶ Customer segmentation
 - ▶ What kinds of customers do I have?
 - ▶ Is a specific customer of a certain kind?
- ▶ Product development
 - ▶ What is a successful product?

CREATIVE ARTIFICIAL INTELLIGENCE (RECIPES, MUSIC, ART, TEXT)

- ▶ e.g. Vondrick, Carl, Hamed Pirsiavash, and Antonio Torralba. “Generating videos with scene dynamics.”

Advances In Neural Information Processing Systems. 2016.⁵

- ▶ Generate an artefact
 - ▶ Generate videos
 - ▶ Generate text
 - ▶ Generate music



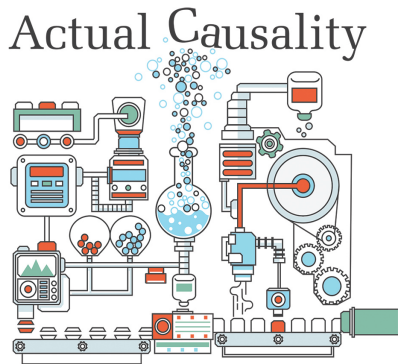
Train Station



⁶http://www.cv-foundation.org/openaccess/content_iccv_2015/papers/Antol_VQA_Visual_Question_ICCV_2015_paper.pdf

CAUSALITY / OFF-POLICY REINFORCEMENT LEARNING

- ▶ Forward Causality
 - ▶ “The effects of causes” - ML Causality
 - ▶ “How much would my sales increase if I send more aggressive e-mails to my customers?”
- ▶ Backward Causality
 - ▶ “The causes of effects” - ML interpretability on causally consistent models
 - ▶ “What was the main cause of the drop in sales?”



Joseph Y. Halpern

REINFORCEMENT LEARNING / BANDITS

- ▶ What if we had access to a simulator?
 - ▶ Or close - we might be able to build one using our data
- ▶ We can search the space of available policies using the simulator and find the best one
- ▶ ... then deploy in the real world



OUR CAPACITY

- ▶ BLG Data Research Center
- ▶ IADS
- ▶ iads@essex.ac.uk, ssamot@essex.ac.uk