

A Short Introduction to Data Science, Big Data Analytics and Statistical Learning

Yu-Tzung Chang and Hsuan-Wei Lee

College of Social Sciences, National Taiwan University

2018.09.13

Data science as a field



Abundance of data

- Volume of data roughly doubles in every three years
 - thank to the technology
 - computers
 - large scale storage devices
 - communication and sensor technologies
- Broad impacts
 - physical, biological, medical sciences
 - all branches of engineering
 - economics, social science, banking and all kinds of commerce
 - government, policy, election
 - environmental science, ecology, outer space
 - sports, recreation

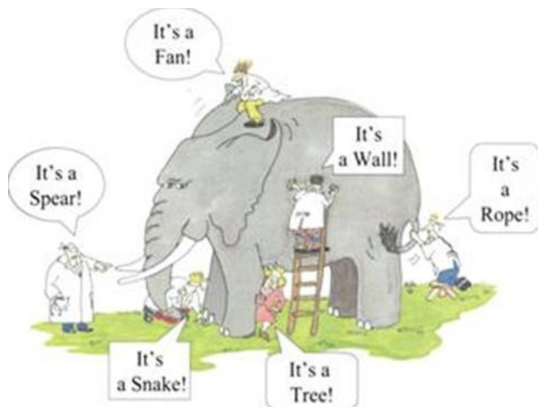
What is big data analytics?

Wikipedia:

Big data is a term used to refer to the study and applications of data sets that are so big and complex that traditional data-processing application software are inadequate to deal with them. Big data challenges include capturing data, data storage, data analysis, search, sharing, transfer, visualization, querying, updating, information privacy and data source.

There are a number of concepts associated with big data: originally there were 3 concepts volume, variety, velocity. Other concepts later attributed with big data are veracity (i.e., how much noise is in the data) and value.

What is data science?



Some quotes from the Internet

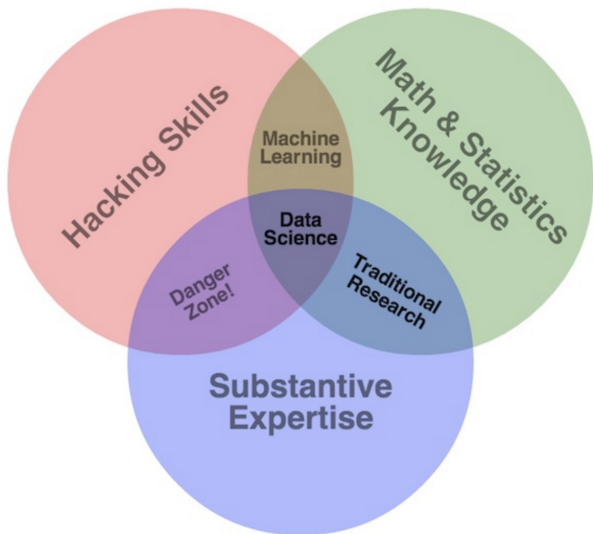
- “A data scientist is a data analyst who lives in California.”
- “A data scientist is a statistician who lives in San Fransisco.”
- “Data science is statistics on a Mac.”
- “A data scientist is someone who is better at statistics than any software engineer and better at software engineering than any statistician.”

Data science and prediction

- “Data science is the study of the generalizable extraction of knowledge from data”*
- A key epistemic requirement for new knowledge is its ability to **predict** and not just **explain**
- Compare: the benefit of knowing things in advance vs. the cost of being wrong

¹Dhar V. Data Science and Prediction, Communications of the ACM, Vol. 56 No.12, December 2013

Data science overview



Related disciplines

- Mathematical and statistical modeling
- Data assimilation
- Estimation theory
- Regression analysis
- Time series analysis
- System identification
- Adaptive optimization
- Inverse problem

Goal of data science – prediction

- Need for prediction is all pervasive
- Understand the patterns of nature (e.g. natural science)
- Prevent damages (e.g. weather prediction, crime prevention, disease propagation)
- Enhance welfare (e.g. medical diagnosis, automation)
- Speculation (e.g. investment, hedge funds, election prediction)

What makes prediction so hard?

- Noise (physical versus social systems)
- Not knowing the right question to ask
- Not having the right/enough data
- Not having enough understanding of the patterns
- Believing in the modeling too much
- Chaotic theory

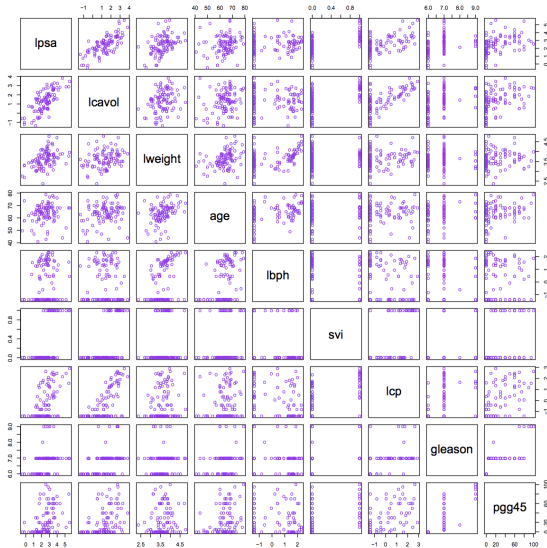
A mathematical and statistical approach

- A very vast discipline
- Key players: data, model, the process of fitting
- Four steps:
 - data mining:
discover basic laws from the analysis of data
 - data assimilation:
an inverse problem to estimate the unknowns
 - model selection:
compare and evaluate models using some criteria
 - generate predictions:
a direct problem to forecast

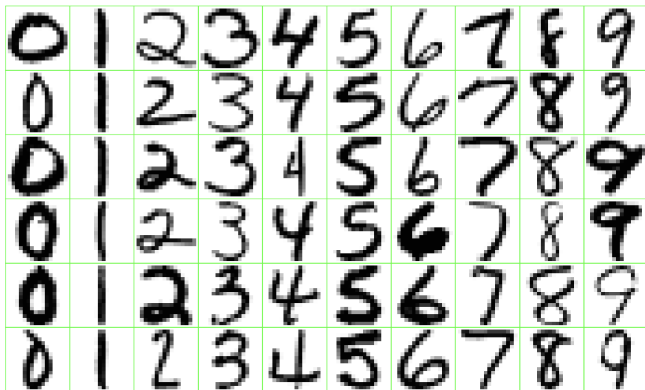
Statistical learning problems (email spam)

	george	you	your	hp	free	hpl	!	our	re	edu	remove
spam	0.00	2.26	1.38	0.02	0.52	0.01	0.51	0.51	0.13	0.01	0.28
email	1.27	1.27	0.44	0.90	0.07	0.43	0.11	0.18	0.42	0.29	0.01

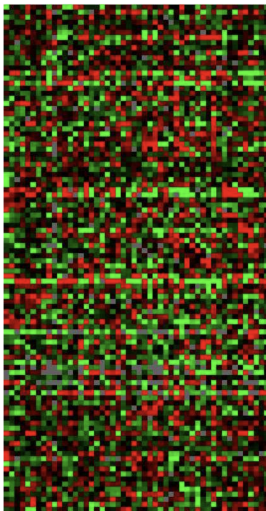
Statistical learning problems (prostate cancer)



Statistical learning problems
(handwritten digit recognition)



Statistical learning problems (DNA expressions microarrays)



SEW399104
SEW398102
SE37701
GNAI
FABP5
SE320394
MAGE17A5E
SE3207172
EST4
SEW377402
HARPM48DNA
SCW369844
EST4
SE471915
MYBPN2D
EST4C61
SE3277401
CHAPL3LME
SE320392
SEW371409
SEW473459
SCW41281
HORMAD98
SCW37068
MTTCCHQND
SE4716
EST4C61
SE320393D
SE488017
SE320391
EST4C61
SE320394
SE320394
PTN1C
SE3703803
SEW370341
SE371068
EST4C61
SE370341
SE377419
SE320392
SEW32030
SEW37064
SEW37034
MAGE17A5E
SE320392
SEW320316
SEW370776
MYO17A1C
WASH190E
SCW37064
EST4C61
SE370394
SE320396
EST4C61
SE488021
SE48836
SEW320315
EST4C61
SEW32036
SE320394
EST4C61
SE320395
SE488140
SE370305
EST4
SEW488740
SMALNLJC
EST4
SEW398311
SEW387187
SE320397
EST4
SE48839
SEW488021
EST4LME
TLPL211LP
SEW48842
SE387079
SEW39832
SEW488271
EST4C61
SEW321025
SCW32030
SEW38162
SE387028
SE377123
SEW38808
EST4C61
SEW321120
SE488317
SE370900
SEW37068
SE37082
SE37084
SE42354

Different language for the same thing

Statistics	Computer Science	Meaning
estimation	learning	using data to estimate an unknown quantity
classification	supervised learning	predicting a discrete Y from $X \in \mathcal{X}$
clustering	unsupervised learning	putting data into groups
data	training sample	$(X_1, Y_1), \dots, (X_n, Y_n)$
covariates	features	the X_i 's
classifier	hypothesis	a map from covariates to outcomes
hypothesis	—	subset of a parameter space Θ
confidence interval	—	interval that contains unknown quantity with a prescribed frequency
directed acyclic graph	Bayes net	multivariate distribution with specified conditional independence relations
Bayesian inference	Bayesian inference	statistical methods for using data to update subjective beliefs
frequentist inference	—	statistical methods for producing point estimates and confidence intervals with guarantees on frequency behavior
large deviation bounds	PAC learning	uniform bounds on probability of errors

¹All of Statistics, Larry A. Wasserman, Springer

Topics of this semester

- Linear regression
- Classification
- Resampling methods
- Linear model selection and regularization
- Moving beyond linearity
- Tree-based methods
- Unsupervised Learning