# **BIG DATA**

(a personal perspective)

Arthur Charpentier

charpentier.arthur@uqam.ca

http://freakonometrics.hypotheses.org/

#### Institut des Actuaires, Paris, May 2014

"Big Data is like teenage sex : everybody talks about it, nobody really knows how to do it, everyone thinks everyone else is doing it, so everyone claims they are doing it..."

Dan Ariely, 2013 facebook.com/dan.ariely/

# **BIG DATA**

#### (an actuarial & a statistical perspective)

#### Arthur Charpentier

charpentier.arthur@uqam.ca

http://freakonometrics.hypotheses.org/

#### Institut des Actuaires, Paris, May 2014

Professor of Actuarial Sciences, Mathematics Department, UQàM
(previously Economics Department, Univ. Rennes 1 & ENSAE Paristech actuary in Hong Kong, IT & Stats FFSA)
PhD in Statistics (KU Leuven), Fellow Institute of Actuaries
MSc in Financial Mathematics (Paris Dauphine) & ENSAE
Editor of the freakonometrics.hypotheses.org's blog
Editor of (forthcoming) Computational Actuarial Science, CRC The R Series

Computational Actuarial Science with R



# Agenda

- Introduction
- Examples of Big Data Issues
- $\circ\,$  Basketball, Maps, Amazon, Netflix, Google, Wikipedia, and the Flu
- Some Thoughts about Big Data
- Defining Big Data
- Volume
- How Big is Big Data?
- $\circ\,$  Correlation, Parametric and Non-Parametric Modeling
- Velocity
- High Frequency Trading
- Variety
- Textming, Graphs and Translation
- Conclusion?

## **Basketball and Optical Tracking Data**

Moneyball 2.0 : How Missile Tracking Cameras Are Remaking The NBA, 2012 fastcodesign.com

Predicting Points and Valuing Decisions in Real Time With NBA Optical Tracking Data. MIT Sloan, Sports Analytics Conference, Feb28-Mar1.

Basketball IRL and *Databall*, 2014 grantland.com



### Basketball and Optical Tracking Data, the Matrix

х	quarte	r game	_clock	ball_	x ball_	y ball_z	off1_ent	off2_ent	off3_ent	off4_ent	off5_ent	off1_x	off2_x	off3_x	off4_x	off5_x	
1	11760	94	4	9.2	4 30.4618	4 15.66440	5.73955	Parker	Neal	Leonard	Bonner	Duncan	31.20810	3.01034	5.71867	22.70688 2	2.66160
2	11760	95	4	9.2	0 30.7489	3 16.56177	6.16081	Parker	Neal	Leonard	Bonner	Duncan	31.55756	2.90998	5.59025	22.60166 2	2.74272
3	11760	96	4	9.1	6 31.1400	3 16.98810	6.08396	Parker	Neal	Leonard	Bonner	Duncan	31.89928	2.81066	5.47863	22.47178 2	2.89309
4	11760	97	4	9.1	2 31.8679	l 17.09691	6.06588	Parker	Neal	Leonard	Bonner	Duncan	32.23151	2.71905	5.39309	22.31212 2	3.04843
5	11760	98	4	9.0	8 31.9917	4 17.04352	5.80229	Parker	Neal	Leonard	Bonner	Duncan	32.54901	2.64163	5.33036	22.14823 2	3.20962
6	11760	99	4	9.0	4 32.5539	6 16.93693	5.61020	Parker	Neal	Leonard	Bonner	Duncan	32.86782	2.56882	5.28070	21.97328 2	3.38824
7	11761	00	4	9.0	0 33.0152	3 16.69451	5.24575	Parker	Neal	Leonard	Bonner	Duncan	33.17220	2.50098	5.24495	21.80377 2	3.58398
8	11761	01	4	8.9	6 33.5985	3 16.52964	4.81964	Parker	Neal	Leonard	Bonner	Duncan	33.47203	2.44498	5.21824	21.61287 2	3.79921
9	11761	02	4	8.9	2 33.7556	) 16.17879	4.53087	Parker	Neal	Leonard	Bonner	Duncan	33.76465	2.39688	5.19696	21.39753 2	4.03398
10	11761	03	4	8.8	8 34.2841	5 15.67568	4.30612	Parker	Neal	Leonard	Bonner	Duncan	34.04971	2.34841	5.17501	21.17353 2	4.28544
11	11761	04	4	8.8	4 34.2679	7 15.54423	3.69690	Parker	Neal	Leonard	Bonner	Duncan	34.32216	2.30455	5.15439	20.93744 2	4.55324
12	11761	05	4	8.8	0 34.6356	15.57366	3.29938	Parker	Neal	Leonard	Bonner	Duncan	34.58345	2.26725	5.14478	20.68329 2	4.83165
	off	1_y o	ff2_y	off3_y	off4_y	off5_y	def1_en	t def2_e	nt def3_e	ent def4_e	ent def5_en	t def1	_x def2	x def3	x def	4_x def5_	x def1_y
1	18.74	699 3.	09636	47.53907	-1.30109	26.34201	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 28.395	75 7.261	6.611	96 19.10	779 18.4907	4 23.89819
2	18.44	052 2.	95225	47.67462	-1.23840	26.21422	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 28.752	253 7.1468	32 6.529	86 19.03	464 18.6538	7 23.47161
3	18.16	281 2.	82953	47.79135	-1.19877	26.05057	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 29.104	03 7.0410	02 6.454	09 18.95	038 18.8215	9 23.04556
4	17.92	252 2.	70648	47.88714	-1.16482	25.87835	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 29.450	96 6.9489	94 6.383	21 18.84	988 18.9927	5 22.61536
	17.70	649 2.	58385	47.96527	-1.12016	25.69984	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 29.770	94 6.8718	88 6.318	62 18.73	036 19.1654	9 22.20077
6	17.52	267 2.	46749	48.03880	-1.06383	25.51339	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 30.086	60 6.813	73 6.265	94 18.60	247 19.3395	2 21.80239
7	17.35	745 2.	36363	48.10369	-1.00097	25.31308	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 30.394	00 6.7620	59 6.219	68 18.45	502 19.5240	1 21.42181
8	17.20	283 2.	29945	48.15793	-0.94003	25.10529	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 30.688	64 6.719	73 6.185	66 18.29	982 19.7108	2 21.06233
9	17.07	512 2.	24979	48,20172	-0.87742	24.88480	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 31.002	25 6.6789	99 6.160	05 18.14	136 19.9062	2 20.72579
10	16.97	912 2.	20291	48.23125	-0.80572	24.66389	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 31.318	60 6.645	73 6.150	86 17.98	013 20.1142	0 20.41285
11	16.90	115 2.	16288	48.25132	-0.72214	24.43810	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 31.630	056 6.6114	47 6.147	76 17.81	079 20.3286	1 20.12016
12	16.84	963 2.	13311	48.26263	-0.63485	24.21707	Livingsto	n Ellingt	on Waite	ers (	Gee Zelle	r 31.931	98 6.5832	26 6.155	69 17.63	143 20.5531	3 19.85651
	def2	_y d	ef3_y	def4_y	def5_y j	prob_pass1	prob_pa	ss2 pro	b_pass3	prob_pas	ss4 prob	pass5	prob_mal	ce pro	ob_miss	prob_T	0 epv
1	7.870	51 36.	60250	2.23678	20.79599	NA	0.000939	464 0.000	5 <u>2</u> 15383 0	0.00078003	394 0.0011 <del>0</del>	80469 2.	06139 <del>6</del> e-(	05 3.530	90 <u>0</u> e-05 4	4.582805e-0	5 0.9694958
2	7.786	49 36.	53292	2.26203	20.82908	NA	0.001121	139 0.000	6538326 0	.00071408	829 0.00092	35663 1.	484419e-0	05 2.593	288e-05 4	4.617669e-0	5 0.9679251
3	7.724	41 36.	46256	2.29366	20.86528	NA	0.001297	789 0.000	7020684 0	.0006851	674 0.00079	77597 1.	005769e-0	05 1.793	289e-05 (	4.619174e-0	5 0.9660304
4	7.677	86 36.	39285	2.33380	20.90569	NA	0.001430	896 0.000	7368078 0	.00069359	951 0.00073	56085 6.	828948e-0	06 1.242	612e-05 (	4.780346e-0	5 0.9639335
5	7.644	93 36.	31915	2.38695	20.95015	NA	0.002003	120 0.000	9443742 0	.00106693	102 0.00079	72303 9.	493768e-0	06 1.762	269e-05 4	4.879118e-0	5 0.9616730
6	7.619	30 36.	24331	2.44617	20.99211	NA	0.002028	606 0.000	9864967 0	.0010422	796 0.00070	23052 6.	369296e-0	06 1.206	092e-05 !	5.125298e-0	5 0.9586997
7	7.605	50 36.	17378	2.51259	21.04188	NA	0.001883	377 0.001	0428576 0	.0009959	774 0.00056	48587 4.	288592e-0	06 8.279	725e-06	5.458791e-0	5 0.9556066
8	7.601	38 36.	10760	2.59764	21.08655	NA	0.001768	323 0.001	1071285 0	.0009526	763 0.00052	94989 2.	934505e-0	06 5.770	166e-06	5.815748e-0	5 0.9525035
9	7.604	26 36.	05122	2.66967	21.11940	NA	0.001577	596 0.001	1333265 0	.00089952	217 0.00054	57132 2.	002034e-0	06 4.016	175e-06	6.117891e-0	5 0.9493675
10	7.616	21 35.	99961	2.74528	21.13542	NA	0.001415	889 0.001	1766095 0	.0008518	788 0.00058	63730 1.	391663e-0	06 2.850	721e-06	6.527485e-0	5 0.9462897
11	7.636	96 35.	94937	2.85717	21.14073	NA	0.001238	335 0.001	1980263 0	.0008016	188 0.00062	05814 9.	694368e-0	07 2.028	000e-06	6.893162e-0	5 0.9432457
12	7.665	33 35.	90460	2.97045	21.14716	NA	0.001071	548 0.001	2006243 0	.00074740	017 0.00063	43324 6.	795564e-0	07 1.451	362e-06 '	7.255664e-0	5 0.9402737

## Collection Data for Maps, and Open Street Map

More than 1 million contributors on openstreetmap.org



See Over 100,000 buildings mapped in Guinea where Ebola broke out mapbox.com and OpenStreetMap mapping progress of Guéckédou to support doctors in Guinea after Ebola outbreak datarep.tumblr.com.



Map by John Snow showing the clusters of cholera cases in the London epidemic of 1854, simonrogers.cartodb.com



Open Data, e.g. Open Street Map

Smart meters could bust marijuana nurseries, baltimoresun.com

## Amazon, Netflix (and Sparse Matrices)



Hastie, 2009 use-R-2009

**Example** : Netflix problem.

We partially observe a matrix of movie ratings (rows) by a number of raters (columns). The goal is to predict the future ratings of these same individuals for movies they have not yet rated (or seen)

Statistical significance of the Netflix challenge arxiv.org

## Amazon, Netflix (and Sparse Matrices)



Hastie, 2009 use-R-2009



## Google, Wikipedia and the Flu

2012, Google answers 100 billion search queries a month.

Detecting influenza epidemics using search engine query data nature.com



Predicting the present with Google Trends people.ischool.berkeley.edu Nowcasting with Google Trends inan emerging market ideas.repec.org

#### PERCENT OF HEALTH VISITS FOR FLU-LIKE SYMPTOMS Mid-Atlantic region

#### Using Google to Monitor the Flu

## Google, Wikipedia and the Flu

Google Flu Trends : The Limits of Big Data bits.blogs.nytimes.com



Wikipedia Usage Estimates Prevalence of Influenza-Like Illness in the U.S in Near Real-Time ploscompbiol.org





## **Big Data in Sciences : Astronomy**



#### Example : Kepler's laws (1609)

1. The orbit of a planet is an ellipse with the Sun at one of the two foci.

2. A line segment joining a planet and the Sun sweeps out equal areas during equal intervals of time (The two shaded sectors A1 and A2 have the same surface area and the time for planet 1 to cover segment A1 is equal to the time to cover segment A2.)

Data = Tycho Brahe's data (687 day intervals)

How Big Data Is Changing Astronomy (Again) theatlantic.com

## **Big Data in Sciences : Astronomy**





More Than Meets the Eye : How the CCD Transformed Science wired.com

"There are two reasons that astronomy is experiencing this accelerating explosion of data. First, we are getting very good at building telescopes that can image enormous portions of the sky. Second, the sensitivity of our detectors is subject to the exponential force of Moore's Law."

How Big Data Is Changing Astronomy (Again) theatlantic.com

### **Big Data in Sciences : Climate**



 $\rightarrow$  Improvement of Climate Model resolution over the four IPCC reports clivar.org

Much more noisy data, need filtering techniques

Ocean Surface Speed in NOAA/GFDL Southern Ocean Simulations





## **Big Data in Sciences : Genomic**

G **A T C** A C A G **A A A** T T C C A G C A T A T G A C A T C C A C G C G C T A G C **C G G** T A T A T G **A A A** T G A G A G G **A T C** A T C A C A C T A T G T G A T G A C A T A C T A G A C **C G G** T G A T G **`ATC**AGGAATTCCAGCATATGAC**ATC** C G C T A G C **C G G** T A T A T G A A G G A T G A G A G C C A C C A C T A T G T G A T G A C A T A C T A G A C **C G G** T ACGATGGATTACAGGAATTCCAGCATATGACA G A G G C C A C G C G C T A G C G C G T A T A T G A A A A G A G G G A C A C C A C T A T G T G A T G A C A T A C T TGACATACACGCGCTAGCCGAGTATATGAGAG ACATGAGAGGGACACCACTATGTGATGAC C C C T A G A C **C G G** T G A T G G A T T A C A G G A A ATGACACCCACGCGCTAGCACG G A A A T G A G A G A G G A A T C C A C 1 ACTAGACCGTTTGTGATGGATTACAGG A A T T C C A G C A T A T G A C A T C C A C **A T C** C T A G C T C CAGGTATATGAAAATGAGAGGGACACCACTATG

#### INDEX

**A A A** OFFSETS: 9, 49, 257, 467, 571 **A T C** OFFSETS: 2, 60, 104, 127, 319, 480, 551 **C G G** OFFSETS: 40, 124, 141, 194, 300, 404 Human Genomes and Big Data Challenges osehra.org

"Fast, efficient genome sequencing machines are spewing out more data than geneticists can analyze"

The DNA Data Deluge spectrum.ieee.org

## **Big Data in Sciences : Demography**



#### US Census, 1880

50,189,209 persons censusrecords.com

8 years to tabulate

#### US Census, 1880



Hollerith Tabulating Machine (ex. IBM)7 weeks to tabulate

# **Big Data History**



More details on the history of Big Data :

- A Very Short History Of Big Data forbes.com
- A Very Short History Of Data Science

#### forbes.com

Big Data and the History of Information Storage winshuttle.com

A brief history of big data, the Noam Chomsky way **cnbc.com** 

## **Big Data and Actuarial Sciences**

Pay-As-You-Drive, and Real-Time Pricing Models

Big Data Is My Copilot, business.time.com



Chaire Stratégie Digitale et Big Data hec.fr

## Going Further

#### Many related topics, such as data viz'...

Where your taxes went this year - and where the cuts were made Public spending by the UK's central government departments, 2011-12



Government spending by department, 2011-12 theguardian.com/datablog/

## **Going Further**

... web/data scraping, open data, etc.

## **Defining Big Data**

"The term itself is vague, but it is getting at something that is real" Jon Kleinberg quoted in newyorker.com

For Viktor Mayer-Schönberger, oii.ox.ac.uk, "Big Data is one where n = all" (no sample, but the entire background population) quoted by Tim Harford, ft.com



# **Defining Big Data**

"Large pools of data that can be brought together and analyzed to discern patterns and make better decisions." McKinsey Report, 2011, mckinsey.com

the 3Vs:

increasing volume (amount of data),
velocity (speed of data in and out),
variety (range of data types and sources)

why not got up to 4 or 5Vs?
veracity? (uncertainty in the data)
value? (making big money weforum.org)



3D Data Management : Controlling Data Volume, Velocity and Variety blogs.gartner.com/doug-laney/

## Data (and Statistics)

Data emerged in 1646 as the plural of the Latin Datum



"Big Data is misnamed in our (academic) world, because data sets have always been big. What is different is that we now have the technology to simply run every scenario." Chris Anderson, 2008 wired.com

"statistics is the grammar of data science. It is crucial to making data speak coherently. But it takes statistics to know whether this difference is significant, or just a random fluctuation. (...) What differentiates data science from statistics is that data science is a holistic approach. We're increasingly finding data in the wild, and data scientists are involved with gathering data, massaging it into a tractable form, making it tell its story, and presenting that story to others." Mike Loukides, 2010 radar.oreilly.com



"In short, the more we learn about biology, the further we find ourselves from a model that can explain it. There is now a better way. Petabytes allow us to say : Correlation is enough. We can stop looking for models. We can analyze the data without hypotheses about what it might show. We can throw the numbers into the biggest computing clusters the world has ever seen and let statistical algorithms find patterns where science cannot." Narinder Singh, 2013 venturebeat.com

see also *I just ran two million regressions*, Xavier Sala-I-Marin, 1997, jstor.org against *We Ran One Regression* David Hendry & Hans-Martin Krolzig, 2004, economics.ouls.ox.ac.uk

Epistemological problem : not more certainty, only high likelihood.



"only an infinite sequence of events (...) could contradict a probability estimate" de Vries, On Probable Grounds filosofie.info

"Data without a model is just noise. But faced with massive data, this approach to science - hypothesize, model, test - is becoming obsolete." Chris Anderson, 2008 wired.com

cf. structural models versus nonparametric statistics, *The Founding of the Econometric Society and Econometrica* jstor.org

$$\begin{cases} C_{t} = \gamma_{10} + \gamma_{11}P_{t} + \gamma_{12}P_{t-1} + \beta_{11}W_{t} + \varepsilon_{1t} \\ I_{t} = \gamma_{20} + \gamma_{21}P_{t} + \gamma_{22}P_{t-1} + \beta_{21}K_{t-1} + \varepsilon_{2t} \\ W_{t} = \gamma_{30} + \gamma_{31}A_{t} + \beta_{31}X_{t} + \beta_{32}X_{t-1} + \varepsilon_{3t} \\ X_{t} = C_{t} + I_{t} + G_{t} \\ P_{t} = X_{t} - T_{t} - W_{t} \\ K_{t} = K_{t-1} + I_{t} \end{cases}$$

#### **Parametric and Non-Parametric Statistics**



 $\rightarrow$  parametric model versus nonparametric model, and machine learning See Breiman's *Statistical Modeling : the Two Cultures* jstor.org

## **Statistics and Data Mining**

"Data mining, more stuffily knowledge discovery in databases, is the art of finding and extracting useful patterns in very large collections of data. It's not quite the same as machine learning, because, while it certainly uses ML techniques, the aim is to directly guide action (praxis!), rather than to develop a technology and theory of induction. In some ways, in fact, it's closer to what statistics calls exploratory data analysis, though with certain advantages and limitations that come from having really big data to explore."

Cosma Shalizi, 2013 vserver1.cscs.lsa.umich.edu/~crshalizi/





Interactive Graphic : How Big Is Big Data ?, 2014 businessweek.com

## Statistics, Significance and *p*-values

"A key issue with applying small-sample statistical inference to large samples is that even minuscule effects can become statistically significant. The increased power leads to a dangerous pitfall as well as to a huge opportunity. The issue is one that statisticians have long been aware of : the p-value problem. Chatfield (1995, p. 70) comments, question is not whether differences are significant (they nearly always are in large samples), but whether they are interesting. Forget statistical significance, what is the practical significance of the results?" Mingfeng Lin, Henry Lucas, Jr. et Galit Shmueli, 2010 galitshmueli.com

"Are there times, I ask, when you just have too much data? When it gets in the way and confuses things? He seems taken aback by this line of questioning. More data is always better, he says." Stephen Baker, the Numerati.

## Volume, How Big could be Big Data?

Consider some logistic model  $Y_i \sim \mathcal{B}(p_i)$  with  $\frac{p_i}{1-p_i} = \exp[\mathbf{X}_i^{\mathsf{T}}\boldsymbol{\beta}]$ , with

$$\uparrow \begin{bmatrix} \vdots & \vdots & \vdots \\ 1 & X_{1,i} & \cdots & X_{k,i} \\ \vdots & \vdots & & \vdots \\ \leftarrow & k+1 & \rightarrow \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$$

where  $\boldsymbol{X}$  is some  $n \times (k+1)$  matrix.

Monte Carlo simulation, with n = 100,000 and k = 100 (but only two  $\beta_j$ 's are not null).

 $\rightarrow$  use of subsampling techniques to estimate several models, n/10 and n/100.

### Volume, How Big could be Big Data?



## Volume, How Big could be Big Data?

Look at ROC curves, instead of  $\hat{\beta}_j$ 's



- average of the 100 regressions on datasets with n/10 and n/100 observations

## **Insurance : Personalization and Customization**

Recall basic results on ratemaking and risk pooling.

No risk classification, identical premium

	Insured	Insurer
Loss	$\mathbb{E}[S]$	$S - \mathbb{E}[S]$
Average Loss	$\mathbb{E}[S]$	0
Variance	0	$\operatorname{Var}[S]$

## **Insurance : Personalization and Customization**

Perfect classification, (ultra) personalized premium

	Insured	Insurer
Loss	$\mathbb{E}[S oldsymbol{\Omega}]$	$S - \mathbb{E}[S \mathbf{\Omega}]$
Average Loss	$\mathbb{E}[S]$	0
Variance	$\operatorname{Var}\left[\mathbb{E}[S \mathbf{\Omega}]\right]$	$\operatorname{Var}\left[S - \mathbb{E}[S \mathbf{\Omega}]\right]$

$$\operatorname{Var}[S] = \underbrace{\mathbb{E}\left[\operatorname{Var}[S|\mathbf{\Omega}]\right]}_{\to \operatorname{insurer}} + \underbrace{\operatorname{Var}\left[\mathbb{E}[S|\mathbf{\Omega}]\right]}_{\to \operatorname{insured}}$$

### **Insurance : Personalization and Customization**

Imperfect classification, personalized premium

	Insured	Insurer
Loss	$\mathbb{E}[S oldsymbol{X}]$	$S - \mathbb{E}[S oldsymbol{X}]$
Average Loss	$\mathbb{E}[S]$	0
Variance	$\operatorname{Var}\!\left[\mathbb{E}[S oldsymbol{X}] ight]$	$\mathbb{E}\Big[\mathrm{Var}[S oldsymbol{X}]\Big]$

$$\operatorname{Var}[S] = \mathbb{E}\left[\operatorname{Var}[S|\boldsymbol{X}]\right] + \operatorname{Var}\left[\mathbb{E}[S|\boldsymbol{X}]\right]$$
$$= \underbrace{\mathbb{E}\left[\operatorname{Var}[S|\boldsymbol{\Omega}]\right]}_{\operatorname{pooling}} + \underbrace{\mathbb{E}\left[\operatorname{Var}\left[\mathbb{E}[S|\boldsymbol{\Omega}]\middle|\boldsymbol{X}\right]\right]}_{\operatorname{solidarity}} + \underbrace{\operatorname{Var}\left[\mathbb{E}[S|\boldsymbol{X}]\right]}_{\rightarrow \operatorname{insured}}.$$

**Velocity : High Frequency Trading** 

Continuous time model,  $dS_t = \mu S_t dt + \sigma S_t dW_t$ 

 $\rightarrow$  discretized time  $dt = \Delta$ , compound return,  $r_t^{(\Delta)} = \log S_{t+\Delta} - \log S_t$ 

realized volatility over period  $[T - h, T], s_t^2 = \sum_{\tau=0}^{h/\Delta} [r_{t-\tau\Delta}^{(\Delta)}]^2$ 

The Wall Street Code : HFT Whisteblower Haim Bodek on Algorithmic Trading nakedcapitalism.com

High Frequency Trading : Threat or Menace? blogs.hbr.org

A healthy side effect of High Frequency Trading? noahpinionblog.blogspot

The problem with high frequency trading **blogs.reuters.com** 





via The High-Frequency Trading Arms Race : Frequent BatchAuctions as a Market Design Response faculty.chicagobooth.edu

## Variety?



### Variety : No More Datawarehouse?

A Relational Model of Data for Large Shared Data Banks, by Edgar Codd, 1970, seas.upenn.edu



A DATABASE SYSTEM

See NoSQL (e.g. nosql-database.org) that provides a mechanism for storage and retrieval of data, modeled in means other than the tabular relations used in relational databases. Used for big data and real-time web applications.

**Remark** NoSQL stands for "Not only SQL"



Classically, databases where structured (tables, relational databases)

Nowadays, most of the world's data is unstructured (text, image, video, voice)

RTs on Twitter



## Variety : Translating Languages

"Google can translate languages without actually knowing them (given equal corpus data, Google can translate Klingon into Farsi as easily as it can translate French into German) Chris Anderson, 2008 wired.com

translate.google.com can translate back and forth between 71 languages

"I have trouble learning languages, and that's precisely the beauty of machine translation : The most important thing is to be good at math and statistics, and to be able to program. (...) So what the system is basically doing (is) correlating existing translations and learning more or less on its own how to do that with billions and billions of words of text (...) In the end, we compute probabilities of translation" Franz Och, quoted in spiegel.de

### Variety : Text Mining and Sentiment Analysis

Automatically Extracting Dialog Models from Conversation Transcripts ieeexplore.ieee.org

Text Mining Medicine the-scientist.com

Text Mining Gun Deaths Data econometricsbysimulation.com from slate.com (crowdsourced) database

Twitter Mood Predicts the Stock Marker arxiv.org

Easy to use some simple Text Analytics to extract intent from Social Media,

my car didn't start this morning? should be late at work

stuck in bed with back pain, again

 $\longrightarrow$  possible to extract some personal information...

(need to detect jokes, sarcasm, amiguity, non-personal information, etc.)

## A slide on Veracity

"Not everything that counts can be counted, and not everything that can be counted counts."

"Since much of the data deluge comes from anonymous and unverified sources, it is necessary to establish and flag the quality of the data before it is included in any ensemble." Arup Dasgupta, geospatialworld.net

How do you feel about online surveys?







Hype Cycles, 2012 gartner.com

"Big data is a vague term for a massive phenomenon that has rapidly become an obsession with entrepreneurs, scientists, governments and the media" ft.com

huge data means increasing risk of *fake discoveries*, a theory-free analysis of mere correlation is inevitably fragile

"many beats of straw look like needles" Trevor Hastie, quoted in nytimes.com



see Magnetic alignment in grazing and resting cattle and deer pnas.org

Big data : The next frontier for innovation, competition, and productivity, 2011 mckinsey.com/insights

Business Insider : Enterprises Aren't Spending Wildly on Big Data But Don't Know If It's Worth It Yet, 2012 businessinsider.com

Five myths about big data, 2013 washingtonpost.com

What Data Can't Do, 2013 nytimes.com

Big Data Boosts Customer Loyalty. No, Really, 2013 forbes.com

... and some more recent posts and articles (March and April)

Future big data analysts will know everything you did today venturebeat.com

The backlash against big data economist.com

Big data and open data : what's what and why does it matter? theguardian.com

Give us back our statistical datawashingtonpost.com

The rise of big data brings tremendous possibilities and frightening perils washingtonpost.com

Finding a great data scientist can feel like searching for Princess Peach. Shes always in another castle. hnews360.com

The Promise and Peril of Big Data aspeninstitute.org

Eight (No, Nine!) Problems With Big Data nytimes.com

What's Up With Big Data Ethics? forbes.com

Is data privacy an out of date concept? nakedsecurity.sophos.com

Google Flu Trends : The Limits of Big Data bits.blogs.nytimes.com

Wikipedia Usage Estimates Prevalence of Influenza-Like Illness in the U.S in Near Real-Time ploscompbiol.org

Big data : are we making a big mistake? ft.com

Data the buzzword vs. data the actual thing noahpinionblog.blogspot.com

Simplifying Data Analysis and Making Sense of Big Data scientificcomputing.com

Google may be a master at data wrangling, but one of its products has been making bogus data-driven predictions newscientist.com

Big Data Doesn't Have to Mean Big Brother recode.net

How Big Is Big Data? businessweek.com

We are moving from an era of private data and public analyses to one of public data and private analyses andrewgelman.com

Can Big Data Help Us Predict The Next Big (Snow) Storm? business2community.com

You Know Who Else Collected Metadata? The Stasi. propublica.org

Data is data, or are they? stancarey.wordpress.com

