Seven Deadly Sins of Big Data



Richard De Veaux Williams College FTC October 5, 2017

The Original Seven Sins

- Lust
- Gluttony
- Greed
- Sloth
- Wrath
- Envy

• Pride



Thanks to Ewan's Corner, Source Unknown

What is Big Data?



How Big Is Big?





What makes Big data Big?





What is it really?

7



What is it really?



What is it really?



So, What are the Seven Deadly Sins of Big Data?

1. Failing to Define the Problem



Production Analysis

- Statistician studies properties for 74000 recent samples from production
 - ♦Viscosity
 - ✦Max temperature to failure
 - ✦Max pressure to failure
 - ◆Density
 - ✦Load at failure
 - ✦Stress at failure









Three Groups!!!





"Do you know we make 3 products?"

Not Fully Understanding the Problem

Ingot cracking

- •3935 30,000 lb. Ingots
- •Up to 25% cracking rate
- •\$30,000 per recast

•90 potential explanatory variables

- •Water composition (reduced)
- Metal composition
- Process variables
- •Other environmental variables









Ingots – First Tree



We know that!! – some alloys are hard to make. That's why we gave you the data in the first place.

Second Tree



What do you think is in those alloys?

One More Time

Looks like Chrome (!?)

Really?

- Did that "solve" the problem?
 - No, but... experimental design
 Enabled us to *focus* on the important variables



Oh, that's funny! -Issac Asimov

Mosaic Plot



By Hour...



2. Underestimating Data Preparation



Data Preparation

60% to 95% of the time is spent preparing the data



87.1% of all statistics are made up

Why is This Hard?

PVA is a philanthropic organization, sanctioned by the US Govt to represent the disabled veterans

They send out 4 million "free gifts", every 6 weeks

- And hope for donations

Data were used for the KDD 1998 cup

- 200,000 donors (100,000 training, 100,000 test)
- 479 (!) demographic variables
 - -Past giving, income, age etc etc etc
- Recent campaign (only for training set)
 - -Did they give? (Target B)

Who should get the current mailing?

- -What's a cost effective strategy
- -How much did they give (Target D)







What's Hard Exactly?

🕷 QuickTi	ime Playe	r File	Edit Vi	ew Windo	w Help			_		_	-			_		16 A	0 *	() ()	90%	E T	ue 4:20 PM	Q :=	5
000										÷ PV	A.jmp												
■ PVA.jmp	* 4		· ODAT	ED OSOUR	TCODE	STAT	ZIP	MAILCO	PVASTA	DOB	NOEX	RECINH	REC P3	RECPG	RECSWE	MDMA	DOMA	CLUSTE	AGE	AGEFL	HOMEOW	CHILD03	
			1 8	01 GRI	0) IL	6108			3712	0					XXXX	T2	36	60				1
			2 9	101 BOA	1	CA	9132			5202	0					XXXX	S1	14	46	E	н		
			3 9	DO1 AMH	1	NC	2701			0	0					XXXX	R2	43			U		
			4 8	701 BRY	C	CA	9595			2801	0					XXXX	R2	44	70	E	U		
			5 8	501	C	FL	3317			2001	0	x	x			XXXX	S2	16	78	E	н		
			6 9	101 CWR	C	AL	3560			0	0					XXXX	T2	40					
			7 8	701 DRK	C	IN	4675			6001	0					XXXX	T2	40	38	E	н		
			8 9	101 NWN	C	LA	7061			0	0					XXXX	T2	39			U		
			9 8	301 LIS	1	IA	5103			0	0					XXXX	R2	45			U		
Columns (491/1)		10 9	101 MSD	1	TN	3712			3211	0					XXXX	T1	35	65	1			
ODATEDW	0		11 9	SO1 AGR	C	KS	6733			0	0					XXXX	R3	53	•		U		
SOURCE			12 9	501 CSM	1	IN	4622			2301	0					XXXX	S2	17	75	E	U		
L TCODE			13 8	01 ENQ	C	MN.	5647			2603	0					XXXX	R3	51	72		н		
STATE			14 9	201 HCC	1.1.1.1	LA	7079			0	0	х				XXXX	T2	40	•				
L ZIP			15 9	301 USB	1	UT	8472			2709	0					XXXX	T1	35	70	E	н		
MAILCODE			16 9	101 FRC	- 1	CA	9005			0	0					XXXX	U1	2			н		
DOB			17 9	101 RKB	C	MI	4806			5401	0					XXXX	S2	20	44	E	U		
NOEXCH			18 8	BO1 PCH	2	2 IL	6237			5201	0					XXXX	R2	43	46	E	U		
& RECINHSE			19 8	SO1 AMB	28	FL	3281	В		3601	0					XXXX	S2	16	62	E	н		
A RECP3			20 9	501 L15	1	NC	2785			0	0					XXXX	C2	27	•				
RECPGVG			21 8	701 BBK	2	MN	5512			3601	0					XXXX	S1	12	62	E	н		
RECSWEEP	-		22 9	501 L21	1	MI	4924			1601	0					XXXX	R2	43	82	E	U		
MDMAUD			23 9	101 SYN	(FL	3398			0	0					XXXX	T2	40					
			24 9	301 L01	2	2 IL	6004			2311	0					XXXX	C1	22	74				
AGE			25 9	501 MOP	C	MN	5504			5201	0					XXXX	T1	35	46	E	н		
			26 9	101 UCA	C	CA	9352			4307	0					XXXX	T1	35	54	1	н		
HOMEOWNR			27 9	601 ESN	C	IL	6009			5601	0					XXXX	S1	13	42	E	н		
CHILD03			28 9	201 L01	1	MO	6446			1401	0					XXXX	T1	35	84	E	н		
L CHILD07			29 9	101 IMP	C	XT (7738			4809	0					XXXX	T1	35	49	E	н	M	٢
Desig			30 9	101 AVN	0	IL	6224			6001	0					XXXX	T2	40	38	E	U		
Hows	~		31 9	001 SYN	C	TX	7754			0	0					XXXX	T1	35					
All rows 94,	,649		32 9	501 RMG	1	MO	6303			2601	0					XXXX	U3	8	72	E	н		
Excluded	0		33 9	501 DNA	28	NC	2710			1401	0					XXXX	C2	25	84	E	U		
Hidden	0		34 9	201 L04	3	FL	3314			2904	0					XXXX	S1	15	69	E	U		
Labelled	0		35 9	101 AML		OR	9700			0	0					XXXX							
			36 9	501		OR	9770			0	0		x			XXXX	B1	42					
			37 9	501 AIR		TX	7874			2901	0					XXXX	S1	12	69	E	н		
			38 8	501 DUR		CA	9027			1002	0					XXXX	S1	11	88		U		
	11		39 9	001 LHJ		MN	5511			2301	0					XXXX	S2	17	75	E	н		
			40 8	01 WKB		MN	5506			1311	0					XXXX	T2	36	84	E	Н		
			41 9	301 AGR			6160			2801	0					XXXX	C2	28	70	1	-		
			12 9	701 STI		MI	4950			0	0					YYYY	111	20					

T-Code

Distributions								
TCODE								
	v Quantil	es		Summary Statistics				
	100.0% 99.5% 97.5% 90.0% 75.0% 50.0% 25.0% 10.0% 2.5% 0.5% 0.0%	maximum quartile median quartile	72002 1002 28 28 2 1 0 0 0 0 0	Mean Std Dev Std Err Mean Upper 95% Mean Lower 95% Mean N	54.413105 957.50024 3.1122959 60.513171 48.313039 94649			

Zoom of Boxplot



Transformation?

Histogram of log10(tcode + 0.01)



Maybe Categories?



What is it, Anyway?

I-Code	litle						
0	-	16	DEAN	48	CORPORAL	109	LIC.
1	MR.	17	JUDGE	50	ELDER	111	SA.
1001	MESSRS.	17002	JUDGE & MRS.	56	MAYOR	114	DA.
1002	MR. & MRS.	18	MAJOR	59002	LIEUTENANT & MRS.	116	SR.
2	MRS.	18002	MAJOR & MRS.	62	LORD	117	SRA.
2002	MESDAMES	19	SENATOR	63	CARDINAL	118	SRTA.
3	MISS	20	GOVERNOR	64	FRIEND	120	YOUR MAJESTY
3003	MISSES	21002	SERGEANT & MRS.	65	FRIENDS	122	HIS HIGHNESS
4	DR.	22002	COLNEL & MRS.	68	ARCHDEACON	123	HER HIGHNESS
4002	DR. & MRS.	24	LIEUTENANT	69	CANON	124	COUNT
4004	DOCTORS	26	MONSIGNOR	70	BISHOP	125	LADY
5	MADAME	27	REVEREND	72002	REVEREND & MRS.	126	PRINCE
6	SERGEANT	28	MS.	73	PASTOR	127	PRINCESS
9	RABBI	28028	MSS.	75	ARCHBISHOP	128	CHIEF
10	PROFESSOR	29	BISHOP	85	SPECIALIST	129	BARON
10002	PROFESSOR & MRS.	31	AMBASSADOR	87	PRIVATE	130	SHEIK
10010	PROFESSORS	31002	AMBASSADOR & MRS.	89	SEAMAN	131	PRINCE AND PRINCESS
11	ADMIRAL	33	CANTOR	90	AIRMAN	132	YOUR IMPERIAL MAJESTY
11002	ADMIRAL & MRS.	36	BROTHER	91	JUSTICE	135	M. ET MME.
12	GENERAL	37	SIR	92	MR. JUSTICE	210	PROF.
12002	GENERAL & MRS.	38	COMMODORE	100	М.		
13	COLONEL	40	FATHER	103	MLLE.		
13002	COLONEL & MRS.	42	SISTER	104	CHANCELLOR		
14	CAPTAIN	43	PRESIDENT	106	REPRESENTATIVE		
14002	CAPTAIN & MRS.	44	MASTER	107	SECRETARY		
15	COMMANDER	46	MOTHER	108	LT. GOVERNOR		
15002	COMMANDER & MRS.	47	CHAPLAIN				

3. Ignoring What's Not There



Depression Clinical Trial Study

- Designed to study antidepressant efficacy
 - Measured via Hamilton Rating Scale
- Side effects
 - Sexual dysfunction
 - Misc safety and tolerability issues
- 428 patients
- Two antidepressants + placebo

The Usual Suspects



What's Missing?



Now Automatic



4. Falling in Love With Your Models



Predicting Malignancy

- Breast cancer data from mammograms
 - Error rates by trained radiologists are near 25% for both false positives and false negatives
- Newer equipment is prohibitively expensive for the developing world
- Early detection of breast cancer is crucial
- Cumulative type I error over a decade is near 100% leading to needless biopsies
- Trees did even worse than radiologists



Combining Models

- Bagging (Bootstrap Aggregation)
 - Bootstrap a data set repeatedly
 - Take many versions of same model (e.g. tree)
 - Random Forest Variation
 - Form a committee of models
 - Take majority rule of predictions
- Boosting
 - Create repeated samples of weighted data
 - Weights based on misclassification
 - Combine by majority rule, or linear combination of predictions

Random Forest Wins

	False Positives	False Negatives			
Simple Tree	32.20%	33.70%			
Neural Network	25.50%	31.70%			
Boosted Trees	24.90%	32.50%			
Boostrap Forest	19.30%	28.80%			
Radiologists	22.40%	35.80%			



5. Using Bad Data



Where Do They Come From?

✦The *Times* of London reports 63.2kg bat

- The mysterious moving decimal
- ♦40% of doctors born on Veteran's Day 1911
 - Data Entry
- ✦Mars Lander lost -- \$125 M
 - Confusion of acceleration units (metric newtons/sec vs. English pound/sec) Everywhere

Ovarian Cancer cure published in Lancet

• Differences due to lab practice, not treatment

Data Quality



Data Quality



Data Quality



What Can They Do?





Study by Large Credit Card Issuer

Entire effect was due to one customer who charged \$3M

6. Confusing Correlation and Causation



Poor George

George Box "All models are wrong, but some are useful"

Peter Norvig

"All models are wrong, and increasingly you can succeed without them."



Really?

Chris Anderson "With enough data, the numbers speak for themselves"

The End of Science

The End of Science

0

The quest for knowledge used to begin with grand theories. Now it begins with massive amounts of data. Welcome to the Petabyte Age.

• Shorter men have greater risk of heart attack



• Who is more likely to have a heart attack in 10 yrs?





• High popsicle sales are strongly correlated with shark





• People who believe in alien abductions prefer Pepsi to Coke



7. Not Taking Your Anti-Hubristines



Data Science

Data Science Venn Diagram v2.0

Data Science

Computer
ScienceMachine
LearningMath and
StatisticsUnicornImage: Computer ScienceImage: Computer ScienceUnicornTraditional
ResearchImage: Computer ScienceSubject Matter ExpertiseSubject Matter Expertise

Copyright © 2014 by Steven Geringer Raleigh, NC. Permission is granted to use, distribute, or modify this image, provided that this copyright notice remains intact



Google Predicts Flu Before CDC!!



Data from Google Inc. Last updated: Aug 19, 2015

©2014 Google - Help - Terms of Service - Privacy - Disclaimer - Discuss

Or... Did They?

BIG DATA

The Parable of Google Flu: Traps in Big Data Analysis

David Lazer,^{1,2*} Ryan Kennedy,^{1,3,4} Gary King,³ Alessandro Vespignani^{5,6,3}

n February 2013, Google Flu Trends (GFT) made headlines but not for a reason that Google executives or the creators of the flu tracking system would have hoped. Nature reported that GFT was predicting more than double the proportion of doctor visits for influenza-like illness (ILI) than the Centers for Disease Control and Prevention (CDC), which bases its estimates on surveillance reports from laboratories across the United States (1, 2). This happened despite the fact that GFT was built to predict CDC reports. Given that GFT is often held up as an exemplary use of big data (3, 4), what lessons can we draw from this error?

The problems we identify are not limited to GFT. Research on whether search or social media can



Large errors in flu prediction were largely avoidable, which offers lessons for the use of big data.

run ever since, with a few changes announced in October 2013 (10, 15).

Although not widely reported until 2013, the new GFT has been persistently overestimating flu prevalence for a much longer time. GFT also missed by a very large margin in the 2011-2012 flu season and has missed high for 100 out of 108 weeks starting with August 2011 (see the graph). These errors are not randomly distributed. For example, last week's errors predict this week's errors (temporal autocorrelation), and the direction and magnitude of error varies with the time of year (seasonality). These patterns mean that GFT overlooks considerable information that could be extracted by traditional statistical methods.

WHAT WE CAN LEARN FROM THE EPIC FAILURE OF GOOGLE FLU TRENDS





Google Flu Trends is no longer good at predicting flu, scientists find

The Seven Virtues

- ✦Define the Problem
- ✦Prepare the Data
 - Use Domain Knowledge
- ✦Be Open to New Methods and Models
- ✦Be Aware of Missing Data
- Ensure Data Quality and Ethical Use of Data
- ◆Use Models, not just Associations
- ✦Work in Teams
 - Acknowledge limits to Big Data Analysis

Thank you!!

