

Statistical Modeling: Bigger and Bigger

David Madigan

Rutgers University
stat.rutgers.edu/~madigan



“in data analysis there is no longer any problem of computation”

- Benzécri, 2005

Logistic Regression

- Linear model for log odds of category membership:

$$\log \frac{p(y=1 | \mathbf{x}_i)}{p(y=-1 | \mathbf{x}_i)} = \sum \beta_j x_{ij} = \boldsymbol{\beta} \mathbf{x}_i$$

Maximum Likelihood Training

- Choose parameters (β_j 's) that maximize probability (likelihood) of class labels (y_i 's) given documents (\mathbf{x}_i 's)

$$L(\boldsymbol{\beta}) = p(\boldsymbol{\beta}|D) = \left(\prod_{i=1}^n \frac{1}{1 + \exp(-\boldsymbol{\beta}^T \mathbf{x}_i y_i)} \right)$$

- Tends to overfit
- Not defined if $d > n$
- Feature selection

Shrinkage Methods

- Avoid combinatorial challenge of feature selection
- L1 shrinkage: regularization + feature selection
- Expanding theoretical understanding
- Large scale
- Empirical performance

Ridge Logistic Regression

Maximum likelihood plus a constraint:

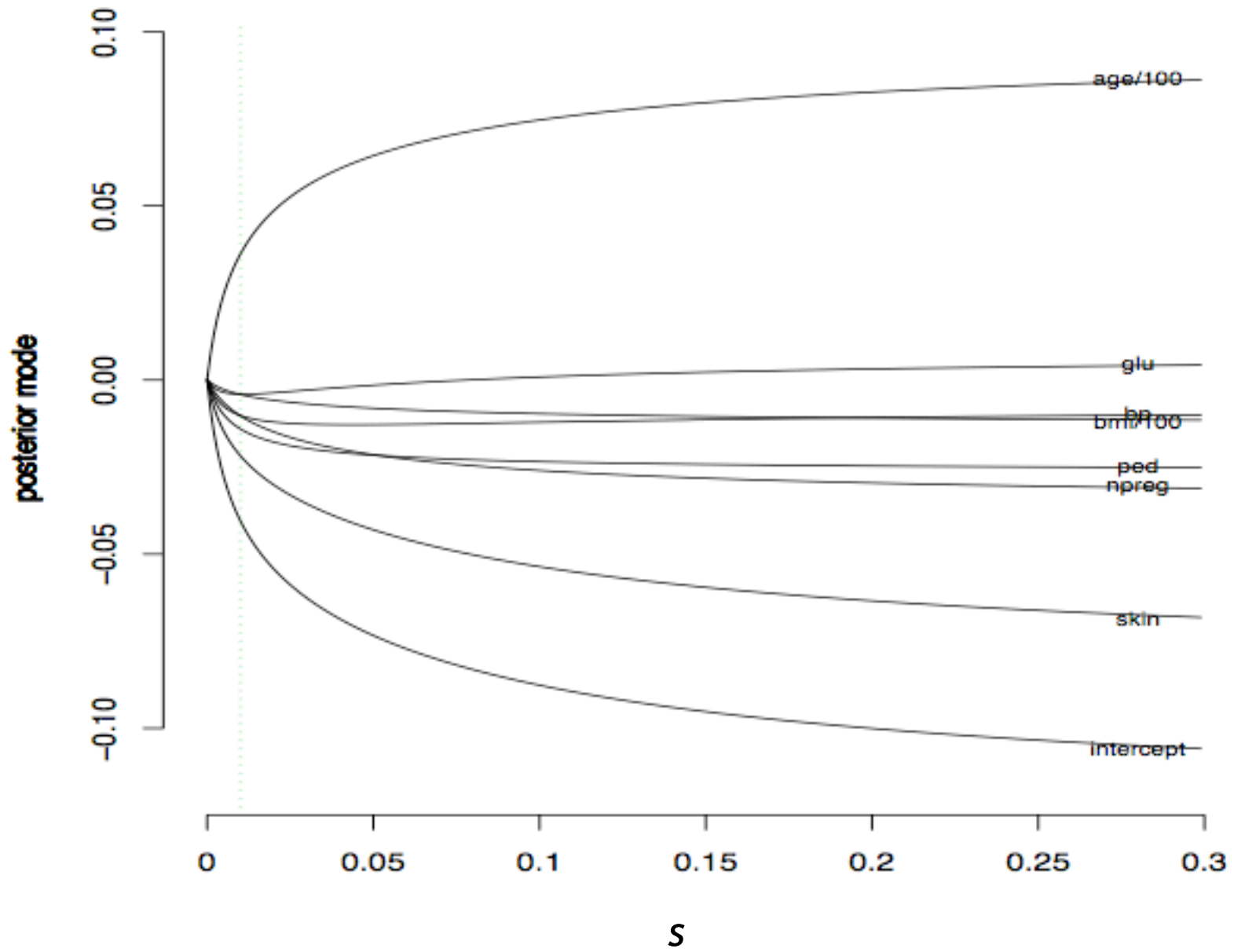
$$\sum_{j=1}^p \beta_j^2 \leq s$$

Lasso Logistic Regression

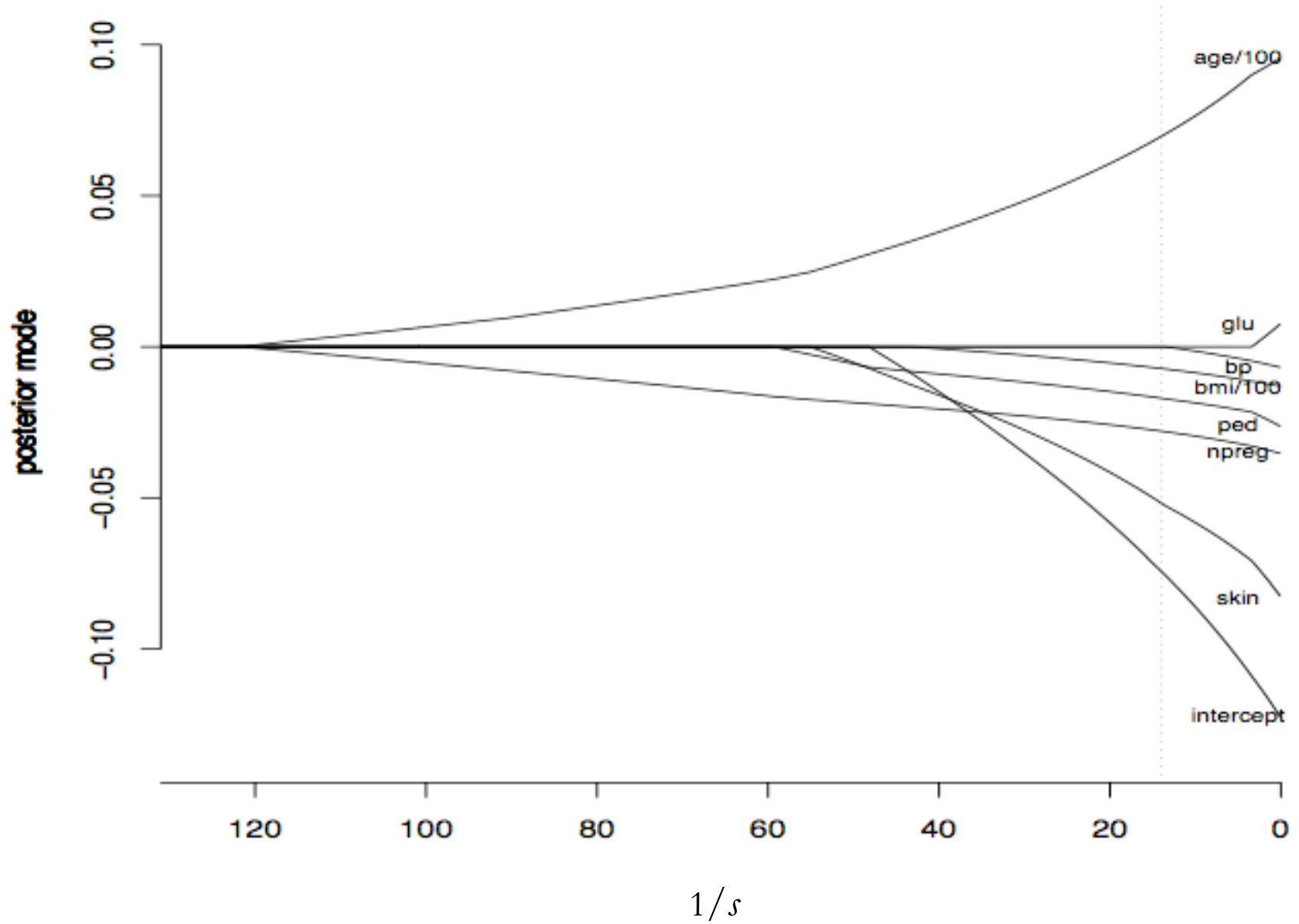
Maximum likelihood plus a constraint:

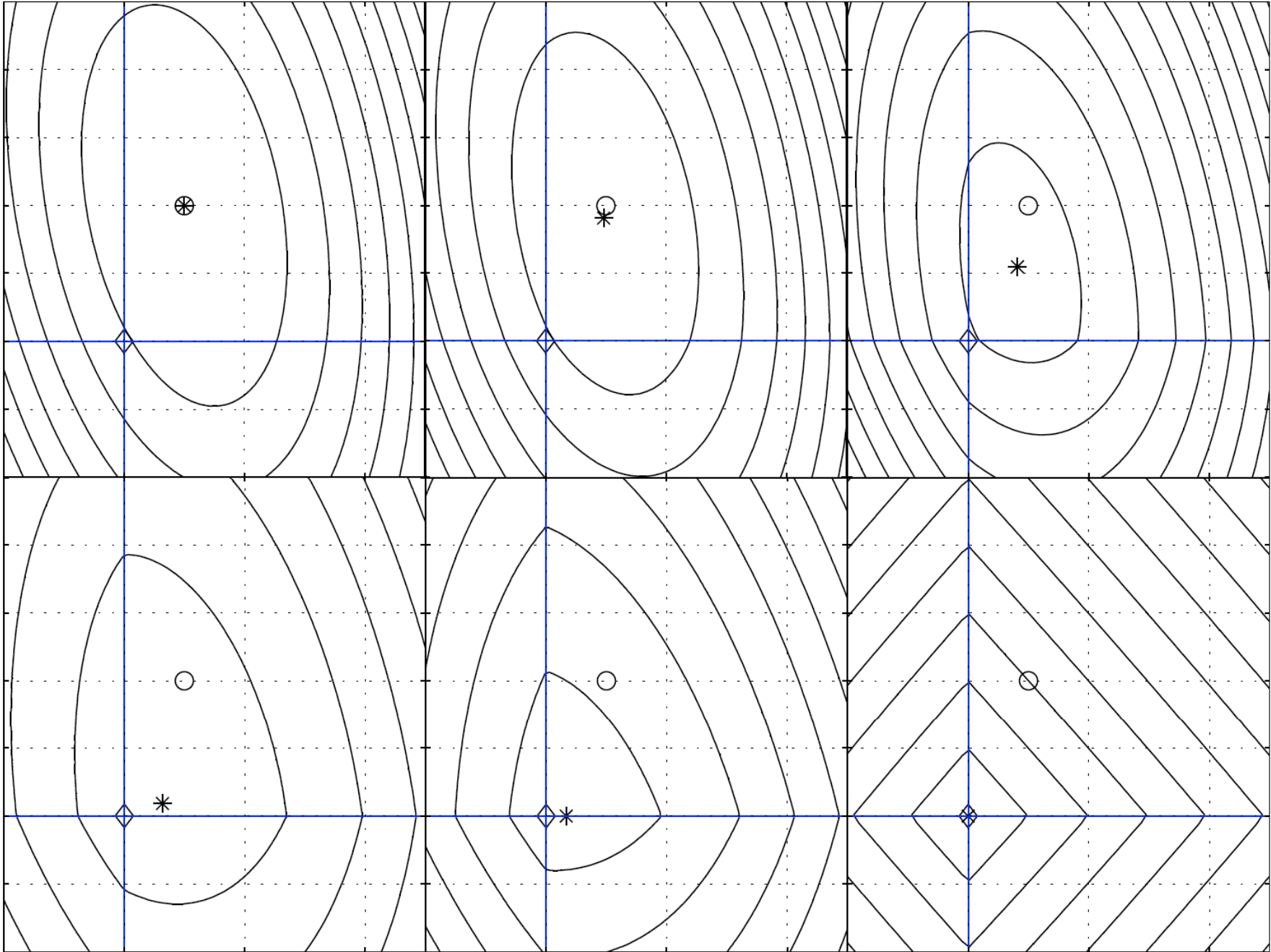
$$\sum_{j=1}^p |\beta_j| \leq s$$

Posterior Modes with Varying Hyperparameter – Gaussian

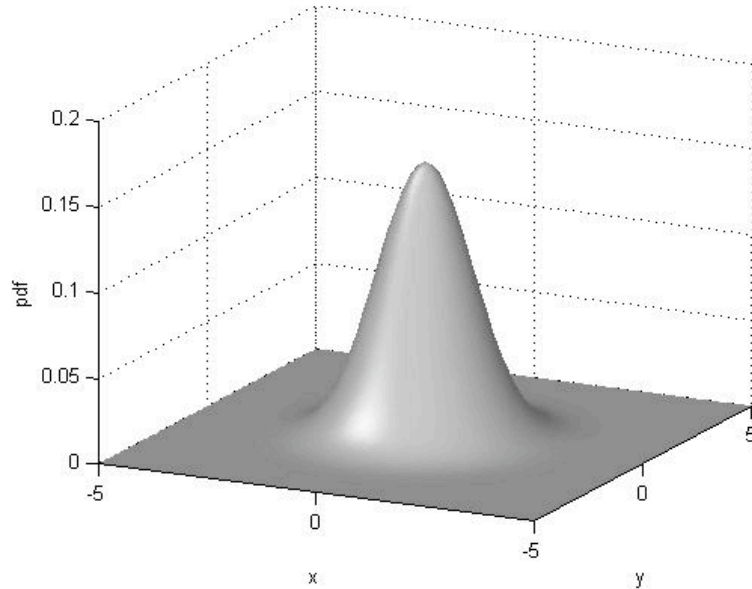


Posterior Modes with Varying Hyperparameter – Laplace

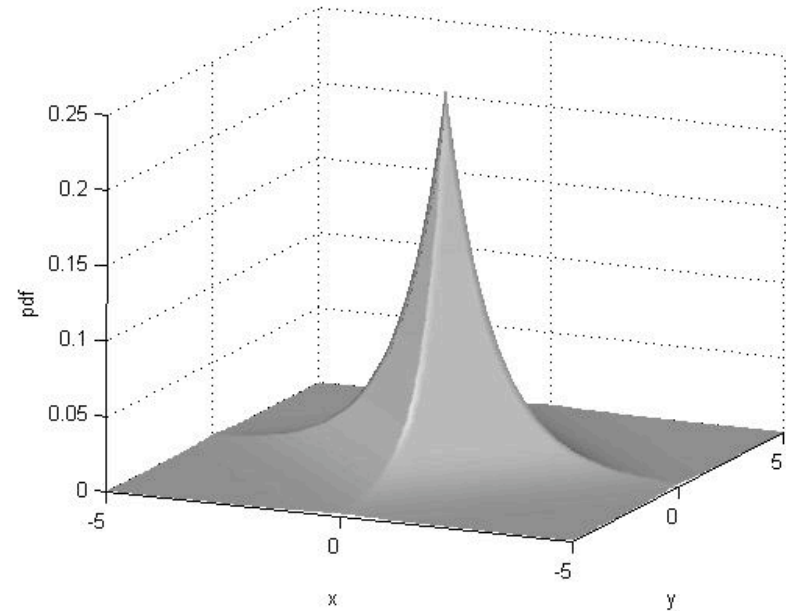




Bayesian Perspective



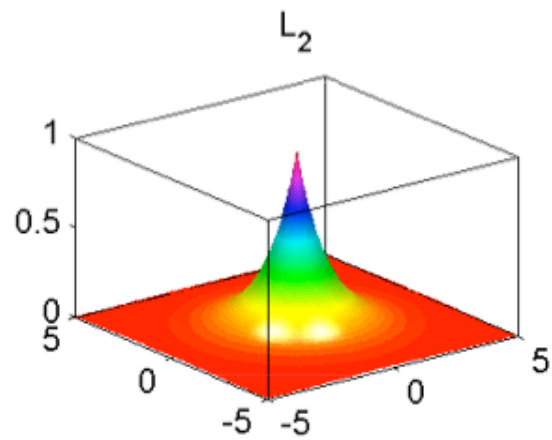
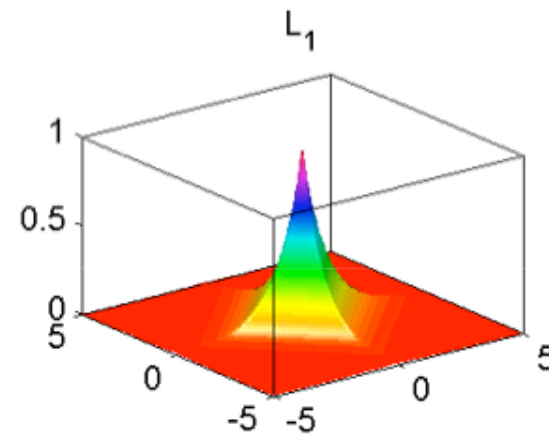
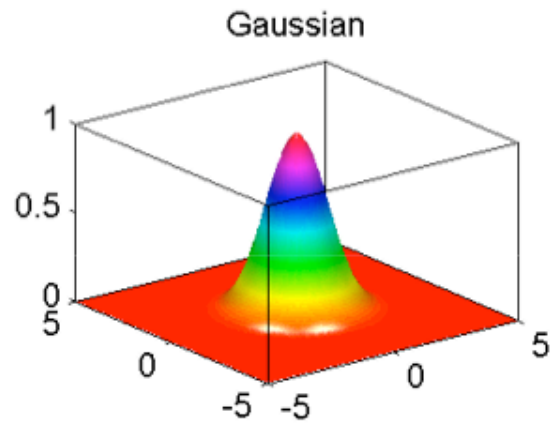
$$\beta_j \sim N(0, \tau^2)$$

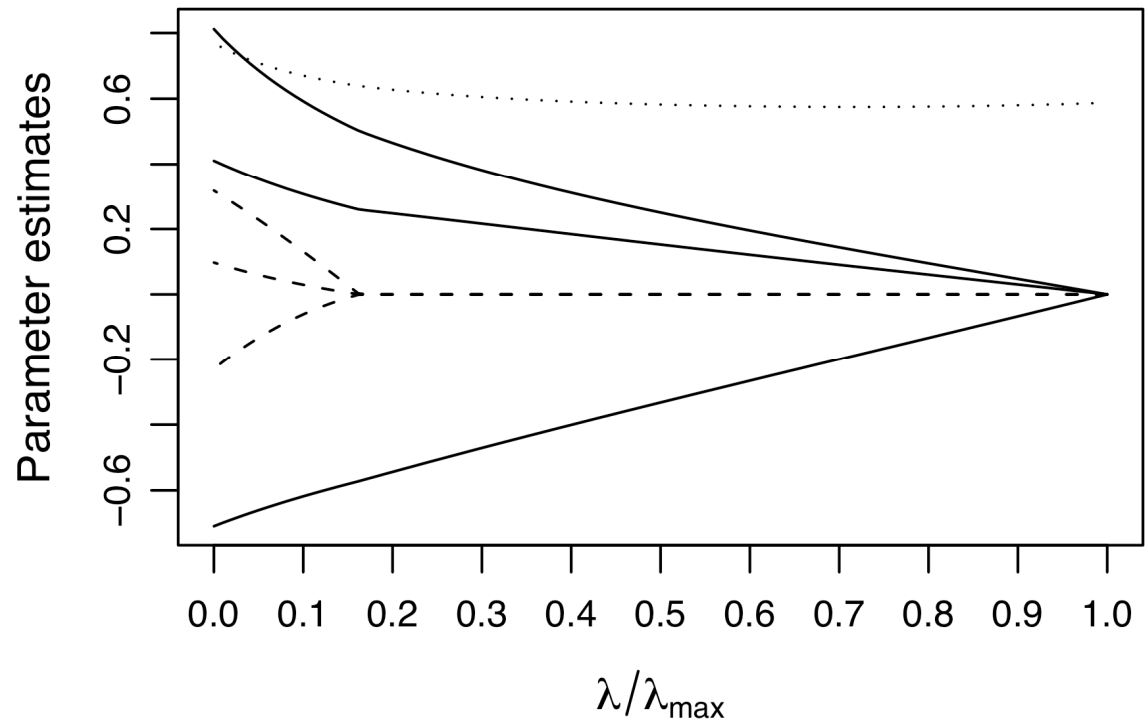


$$\beta_j \sim N(0, \tau_j^2)$$

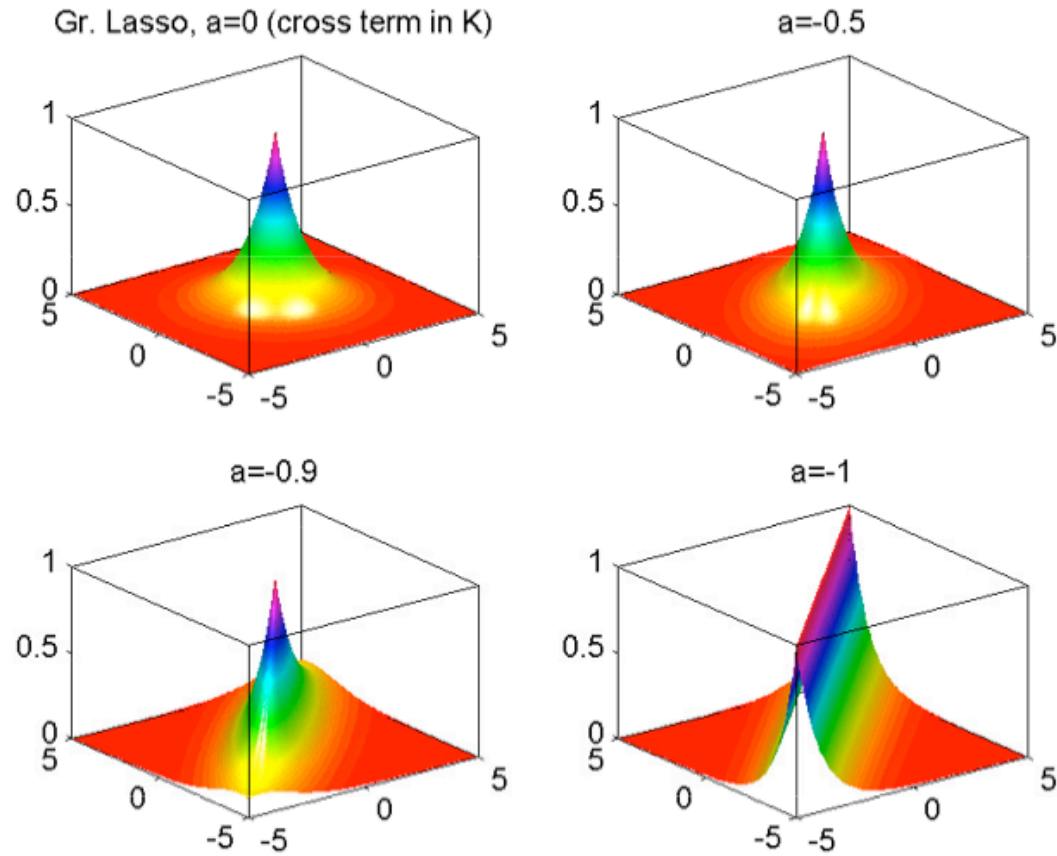
$$\tau_j^2 \sim \exp(\gamma)$$

Group Lasso





“soft fusion”



“Consistency”

- Lasso not always consistent for variable selection
- SCAD (Fan and Li, 2001, JASA) consistent but non-convex
- relaxed lasso (Meinshausen and Bühlmann), adaptive lasso (Wang et al) have certain consistency results
- Zhao and Yu (2006) “irrepresentable condition”

Implementation

- Open source C++ implementation. Compiled versions for Linux, Windows, and Mac (soon)
- Binary and multiclass, hierarchical, informative priors
- Gauss-Seidel co-ordinate descent algorithm
- Fast? (parallel?)
- <http://stat.rutgers.edu/~madigan/BBR>

Aleks Jakulin's results

domain	log-loss / instance						
	BMR	DOT	NB	TAN	MAP	BKT	BK3
krkp	0.09	0.10	-0.29	0.19	<u>0.06</u>	0.11	0.05
monk2	0.65	0.64	-0.65	0.63	<u>0.45</u>	0.60	0.45
tic-tac-toe	<u>0.09</u>	<u>0.08</u>	-0.55	0.49	<u>0.08</u>	0.52	0.07
titanic	0.50	-0.53	0.52	<u>0.48</u>	<u>0.48</u>	<u>0.48</u>	0.48
lenses	0.61	0.72	<u>2.44</u>	-2.99	0.34	0.40	0.40
monk1	-0.50	0.49	0.50	0.09	0.01	0.08	<u>0.02</u>
mushroom	0.00	0.00	-0.01	0.00	0.00	0.00	0.00
shuttle	0.09	0.10	-0.16	0.06	<u>0.07</u>	<u>0.07</u>	<u>0.07</u>
soy-small*	0.27	-0.31	<u>0.00</u>	0.00	0.00	0.00	0.00
wine	<u>0.10</u>	<u>0.09</u>	0.06	<u>0.29</u>	<u>0.19</u>	<u>0.11</u>	<u>0.11</u>
yeast-class*	0.06	0.06	0.01	<u>0.03</u>	-0.25	0.12	0.12
anneal	0.07	0.05	<u>0.07</u>	-0.17	<u>0.11</u>	<u>0.11</u>	<u>0.11</u>
balance-scale	<u>0.20</u>	0.17	0.51	-1.13	0.51	0.51	0.51
lung-cancer*	<u>1.11</u>	1.02	5.41	-6.92	<u>2.37</u>	<u>1.18</u>	<u>1.18</u>
monk3	<u>0.11</u>	0.11	-0.20	<u>0.11</u>	<u>0.11</u>	<u>0.11</u>	<u>0.11</u>
post-op	<u>0.67</u>	0.61	<u>0.93</u>	1.78	<u>0.79</u>	<u>0.67</u>	<u>0.67</u>
promoters*	<u>0.24</u>	0.23	<u>0.60</u>	-3.14	<u>0.59</u>	<u>0.52</u>	<u>0.52</u>
adult	0.28	0.29	-0.42	0.33	0.30	0.30	0.30
audiology*	1.04	1.31	3.55	-5.56	2.24	2.21	2.21
australian	0.33	<u>0.36</u>	0.46	-0.94	<u>0.41</u>	<u>0.37</u>	<u>0.37</u>
breast-LJ	0.55	<u>0.59</u>	<u>0.62</u>	<u>0.89</u>	<u>0.67</u>	<u>0.58</u>	<u>0.58</u>
breast-wisc	0.10	<u>0.12</u>	<u>0.21</u>	<u>0.23</u>	0.21	0.16	0.16
bupa	0.60	<u>0.60</u>	<u>0.62</u>	<u>0.60</u>	<u>0.62</u>	<u>0.61</u>	<u>0.61</u>
car	0.18	<u>0.18</u>	-0.32	<u>0.18</u>	<u>0.19</u>	<u>0.19</u>	<u>0.19</u>
cmc	0.91	0.96	1.00	-1.03	<u>0.93</u>	<u>0.92</u>	<u>0.92</u>
crx	0.33	<u>0.34</u>	<u>0.49</u>	-0.93	<u>0.37</u>	<u>0.35</u>	<u>0.35</u>
ecoli	0.45	0.55	<u>0.89</u>	-0.94	0.85	0.81	0.81
german	0.50	<u>0.51</u>	<u>0.54</u>	-1.04	0.65	<u>0.58</u>	<u>0.59</u>
glass	0.74	<u>0.78</u>	1.25	-1.76	<u>1.12</u>	<u>0.99</u>	<u>0.99</u>
hayes-roth	0.29	<u>0.35</u>	0.46	-1.18	0.45	0.45	0.45
heart	1.01	<u>1.03</u>	1.25	-1.53	1.11	1.09	1.09
hepatitis	0.36	<u>0.39</u>	<u>0.78</u>	-1.31	<u>0.48</u>	<u>0.39</u>	<u>0.39</u>
horse-colic	0.71	<u>0.71</u>	1.67	-5.97	<u>0.83</u>	<u>0.82</u>	<u>0.82</u>
ionosphere	0.19	<u>0.26</u>	0.64	-0.74	0.39	<u>0.30</u>	<u>0.30</u>
iris	0.16	0.24	<u>0.27</u>	<u>0.32</u>	<u>0.27</u>	<u>0.18</u>	<u>0.18</u>
lymph	0.50	<u>0.56</u>	1.10	-1.25	<u>0.98</u>	<u>0.79</u>	<u>0.79</u>
o-ring	0.66	<u>0.80</u>	<u>0.83</u>	<u>0.76</u>	1.41	<u>0.67</u>	<u>0.67</u>
p-tumor*	1.82	1.93	3.17	-4.76	2.65	2.55	2.55
pima	0.46	<u>0.48</u>	<u>0.50</u>	<u>0.49</u>	<u>0.51</u>	<u>0.48</u>	<u>0.48</u>
segment	0.13	<u>0.14</u>	0.38	-1.06	<u>0.17</u>	<u>0.17</u>	<u>0.17</u>
soy-large*	0.25	0.46	0.57	<u>0.47</u>	-0.68	0.66	0.66
spam	0.15	0.16	-0.53	0.32	0.19	0.19	0.19
vehicle	0.54	<u>0.56</u>	-1.78	1.14	0.69	0.66	0.66
voting	0.11	<u>0.13</u>	-0.60	0.53	<u>0.21</u>	<u>0.14</u>	<u>0.14</u>
wdbc	0.09	<u>0.10</u>	0.26	-0.29	<u>0.15</u>	<u>0.13</u>	<u>0.13</u>
zoo*	0.35	-0.47	<u>0.38</u>	<u>0.46</u>	<u>0.40</u>	<u>0.38</u>	<u>0.38</u>
avg rank	2.13	2.87	5.62	5.60	4.74	3.68	3.36

1-of-K Sample Results: brittany-l

Feature Set	% errors	Number of Features
“Argamon” function words, raw tf	74.8	380
POS	75.1	44
1suff	64.2	121
1suff*POS	50.9	554
2suff	40.6	1849
2suff*POS	34.9	3655
3suff	28.7	8676
3suff*POS	27.9	12976
3suff+POS+3suff*POS+Argamon	27.6	22057
All words	23.9	52492

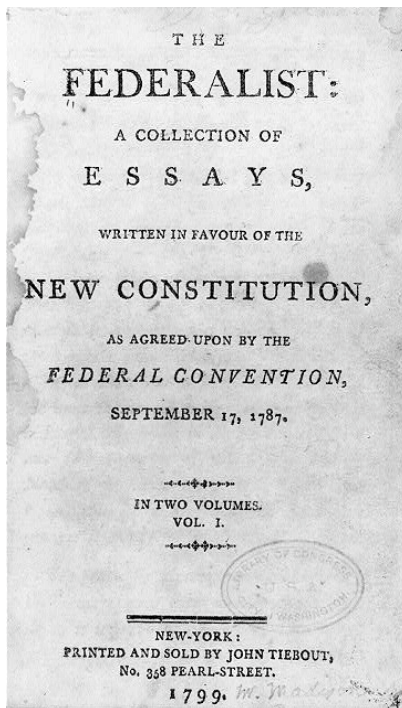
4.6 million parameters

89 authors with at least 50 postings. 10,076 training documents, 3,322 test documents.

BMR-Laplace classification, default hyperparameter

The Federalist

- “The authorship of certain numbers of the ‘Federalist’ has fairly reached the dignity of a well-established historical controversy.” (Henry Cabot Lodge, 1886)
- Historical evidence is muddled



Paper Number	Author
1	Hamilton
2-5	Jay
6-9	Hamilton
10	Madison
11-13	Hamilton
14	Madison
15-17	Hamilton
18-20	Joint: Hamilton and Madison
21-36	Hamilton
37-48	Madison
49-58	Disputed
59-61	Hamilton
62-63	Disputed
64	Jay
65-85	Hamilton



JOURNAL OF THE AMERICAN STATISTICAL ASSOCIATION

Number 302

JUNE, 1963

Volume 58

INFERENCE IN AN AUTHORSHIP PROBLEM^{1,2}

A comparative study of discrimination methods applied
to the authorship of the disputed *Federalist* papers

FREDERICK MOSTELLER

Harvard University

and

Center for Advanced Study in the Behavioral Sciences

AND

DAVID L. WALLACE

University of Chicago

- Used function words with Naïve Bayes with Poisson and Negative Binomial model
- Out-of-sample predictive performance

F. Summing up

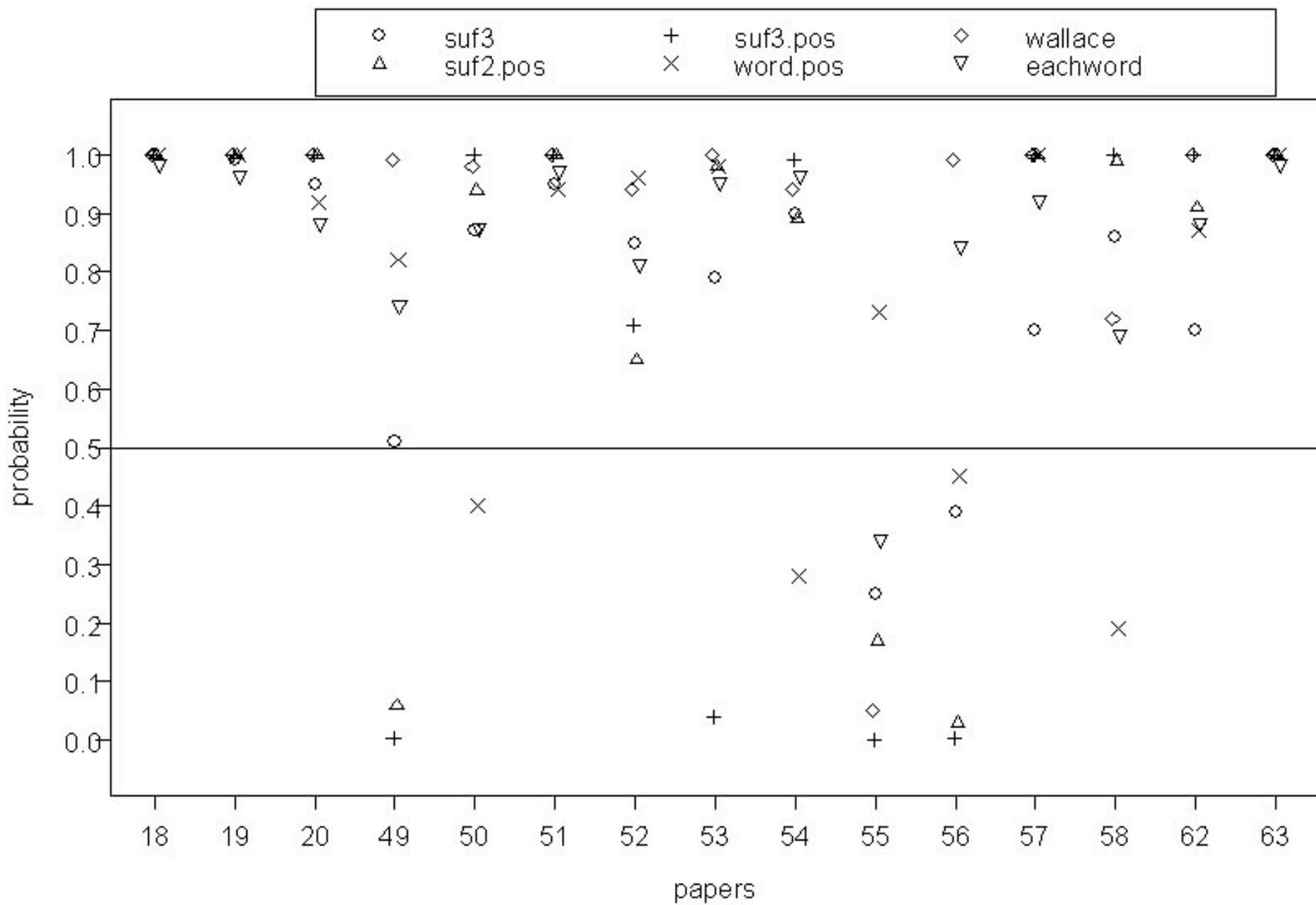
In summary, the following points are clear:

1) Madison is the principal author. These data make it possible to say far more than ever before that the odds are enormously high that Madison wrote the 12 disputed papers. Weakest support is given for No. 55. Support for Nos. 62 and 63, most in doubt by current historians, is tremendous.

Feature Set	10-fold Error Rate
Charcount	0.21
POS	0.19
Suffix2	0.12
Suffix3	0.09
Words	0.10
Charcount+POS	0.12
Suffix2+POS	0.08
Suffix3+POS	0.04
Words+POS	0.08
484 features	0.05
Wallace features	0.05
Words (≥ 2)	0.05
Each Word	0.05

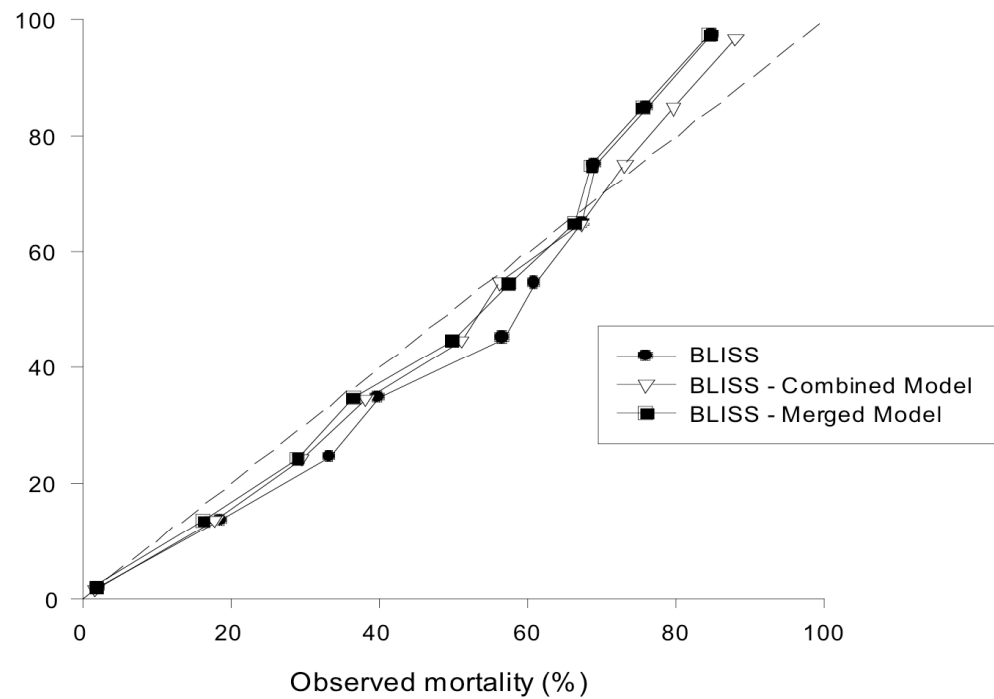
best



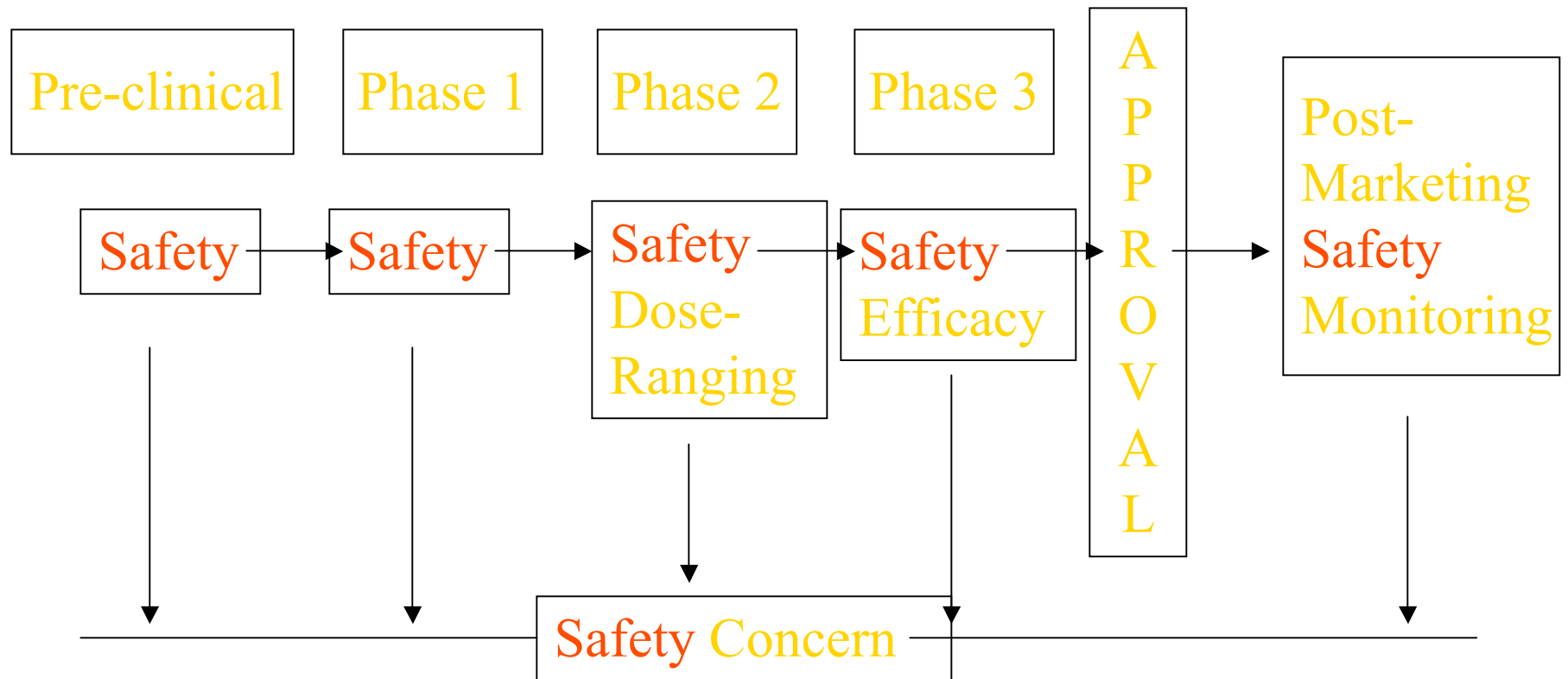


Risk Severity Score for Trauma

- Standard “ICISS” score poorly calibrated
- Lasso logistic regression with 2.5M predictors:



Safety in Lifecycle of a Drug/Biologic product



Databases of Spontaneous ADRs

- FDA Adverse Event Reporting System (AERS)
 - Online 1997 – replace the SRS
 - Over 250,000 ADRs reports annually
 - 15,000 drugs - 16,000 ADRs
- CDC/FDA Vaccine Adverse Events (VAERS)
 - Initiated in 1990
 - 12,000 reports per year
 - 50 vaccines and 700 adverse events
- Other SRS
 - WHO - international pharmacovigilance program

MEDWATCH

For VOLUNTARY reporting of
adverse events, product problems and
product use errors

The FDA Safety Information and
Adverse Event Reporting Program

Page ____ of ____

FDA USE ONLY	
Triage unit sequence #	

A. PATIENT INFORMATION			
1. Patient Identifier	2. Age at Time of Event, or Date of Birth:	3. Sex <input type="checkbox"/> Female <input type="checkbox"/> Male	4. Weight _____ lb or _____ kg
In confidence			

B. ADVERSE EVENT, PRODUCT PROBLEM OR ERROR	
Check all that apply:	
<input type="checkbox"/> Adverse Event	<input type="checkbox"/> Product Problem (e.g., defects/malfunctions)
<input type="checkbox"/> Product Use Error	<input type="checkbox"/> Problem with Different Manufacturer of Same Medicine
2. Outcomes Attributed to Adverse Event (Check all that apply)	
<input type="checkbox"/> Death: _____ (mm/dd/yyyy)	<input type="checkbox"/> Disability or Permanent Damage
<input type="checkbox"/> Life-threatening	<input type="checkbox"/> Congenital Anomaly/Birth Defect
<input type="checkbox"/> Hospitalization - initial or prolonged	<input type="checkbox"/> Other Serious (Important Medical Events)
<input type="checkbox"/> Required intervention to Prevent Permanent Impairment/Damage (Devices)	
3. Date of Event (mm/dd/yyyy)	4. Date of this Report (mm/dd/yyyy)

5. Describe Event, Problem or Product Use Error
6. Relevant Tests/Laboratory Data, including Dates
7. Other Relevant History, including Preexisting Medical Conditions (e.g., allergies, race, pregnancy, smoking and alcohol use, liver/kidney problems, etc.)

C. PRODUCT AVAILABILITY	
Product Available for Evaluation? (Do not send product to FDA)	
<input type="checkbox"/> Yes	<input type="checkbox"/> No
<input type="checkbox"/> Returned to Manufacturer on: _____ (mm/dd/yyyy)	

D. SUSPECT PRODUCT(S)		
1. Name, Strength, Manufacturer (from product label)		
#1		
#2		
2. Dose or Amount Frequency Route		
#1		
#2		
3. Dates of Use (If unknown, give duration) from/to (or best estimate)		5. Event Abated After Use Stopped or Dose Reduced?
#1		#1 <input type="checkbox"/> Yes <input type="checkbox"/> No <input type="checkbox"/> Doesn't Apply
#2		#2 <input type="checkbox"/> Yes <input type="checkbox"/> No <input type="checkbox"/> Doesn't Apply
4. Diagnosis or Reason for Use (Indication)		8. Event Reappeared After Reintroduction?
#1		#1 <input type="checkbox"/> Yes <input type="checkbox"/> No <input type="checkbox"/> Doesn't Apply
#2		#2 <input type="checkbox"/> Yes <input type="checkbox"/> No <input type="checkbox"/> Doesn't Apply
6. Lot #	7. Expiration Date	9. NDC # or Unique ID
#1	#1	
#2	#2	

E. SUSPECT MEDICAL DEVICE		
1. Brand Name		
2. Common Device Name		
3. Manufacturer Name, City and State		
4. Model #	Lot #	5. Operator of Device
Catalog #	Expiration Date (mm/dd/yyyy)	<input type="checkbox"/> Health Professional
Serial #	Other #	<input type="checkbox"/> Lay User/Patient
		<input type="checkbox"/> Other: _____
6. If Implanted, Give Date (mm/dd/yyyy)	7. If Explanted, Give Date (mm/dd/yyyy)	
8. Is this a Single-use Device that was Reprocessed and Reused on a Patient?		
<input type="checkbox"/> Yes <input type="checkbox"/> No		
9. If Yes to Item No. 8, Enter Name and Address of Reprocessor		

F. OTHER (CONCOMITANT) MEDICAL PRODUCTS
Product names and therapy dates (exclude treatment of event)

G. REPORTER (See confidentiality section on back)		
1. Name and Address		
Phone # E-mail		
2. Health Professional? <input type="checkbox"/> Yes <input type="checkbox"/> No	3. Occupation	4. Also Reported to: <input type="checkbox"/> Manufacturer <input type="checkbox"/> User Facility <input type="checkbox"/> Distributor/Importer
5. If you do NOT want your identity disclosed to the manufacturer, place an "X" in this box: <input type="checkbox"/>		

PLEASE TYPE OR USE BLACK INK

Weakness of SRS Data

- Passive surveillance
 - Underreporting
- Lack of accurate “denominator”, only “numerator”
 - “Numerator”: No. of reports of suspected reaction
 - “Denominator”: No. of doses of administered drug
- No certainty that a reported reaction was causal
- Missing, inaccurate or duplicated data

Existing Methods

- Multi-item Gamma Poisson Shrinker (MGPS)
 - US Food and Drug Administration (FDA)
- Bayesian Confidence Propagation Neural Network
 - WHO Uppsala Monitoring Centre (UMC)
- Proportional Reporting Ratio (PRR and aPRR)
 - UK Medicines Control Agency (MCA)
- Reporting Odds Ratios and Incidence Rate Ratios
 - Other national spontaneous reporting centers and drug safety research units

Existing Methods (Cont'd)

- Focus on 2X2 contingency table projections

	<i>AE j = Yes</i>	<i>AE j = No</i>	<i>Total</i>
<i>Drug i = Yes</i>	<i>a=20</i>	<i>b=100</i>	<i>120</i>
<i>Drug i = No</i>	<i>c=100</i>	<i>d=980</i>	<i>1080</i>
<i>Total</i>	<i>120</i>	<i>1080</i>	<i>1200</i>

- 15,000 drugs * 16,000 AEs = 240 million tables
- Most $N_{ij} = 0$, even though $N_{.}$ very large

The Different Measures

<i>Measure of Association</i>	<i>Formula</i>	<i>Probabilistic Interpretation</i>
RR Relative Risk*	$\frac{a * (a + b + c + d)}{(a + c) * (a + b)}$	$\frac{\Pr(ae drug)}{\Pr(ae)}$
PRR Proportional Reporting Ratio	$\frac{a / (a + b)}{c / (c + d)}$	$\frac{\Pr(ae drug)}{\Pr(ae \neg drug)}$
ROR Reporting Odds Ratio	$\frac{a / c}{b / d}$	$\frac{\Pr(ae drug) / \Pr(\neg ae drug)}{\Pr(ae \neg drug) / \Pr(\neg ae drug)}$
Information Component	$\text{Log}_2 \frac{a * (a + b + c + d)}{(a + c) * (a + b)}$	$\log_2 \frac{\Pr(ae drug)}{\Pr(ae)}$

Relative Reporting Ratio

$$(RR_{ij} = N_{ij} / E_{ij})$$

- Advantages

- Simple
- Easy to interpret

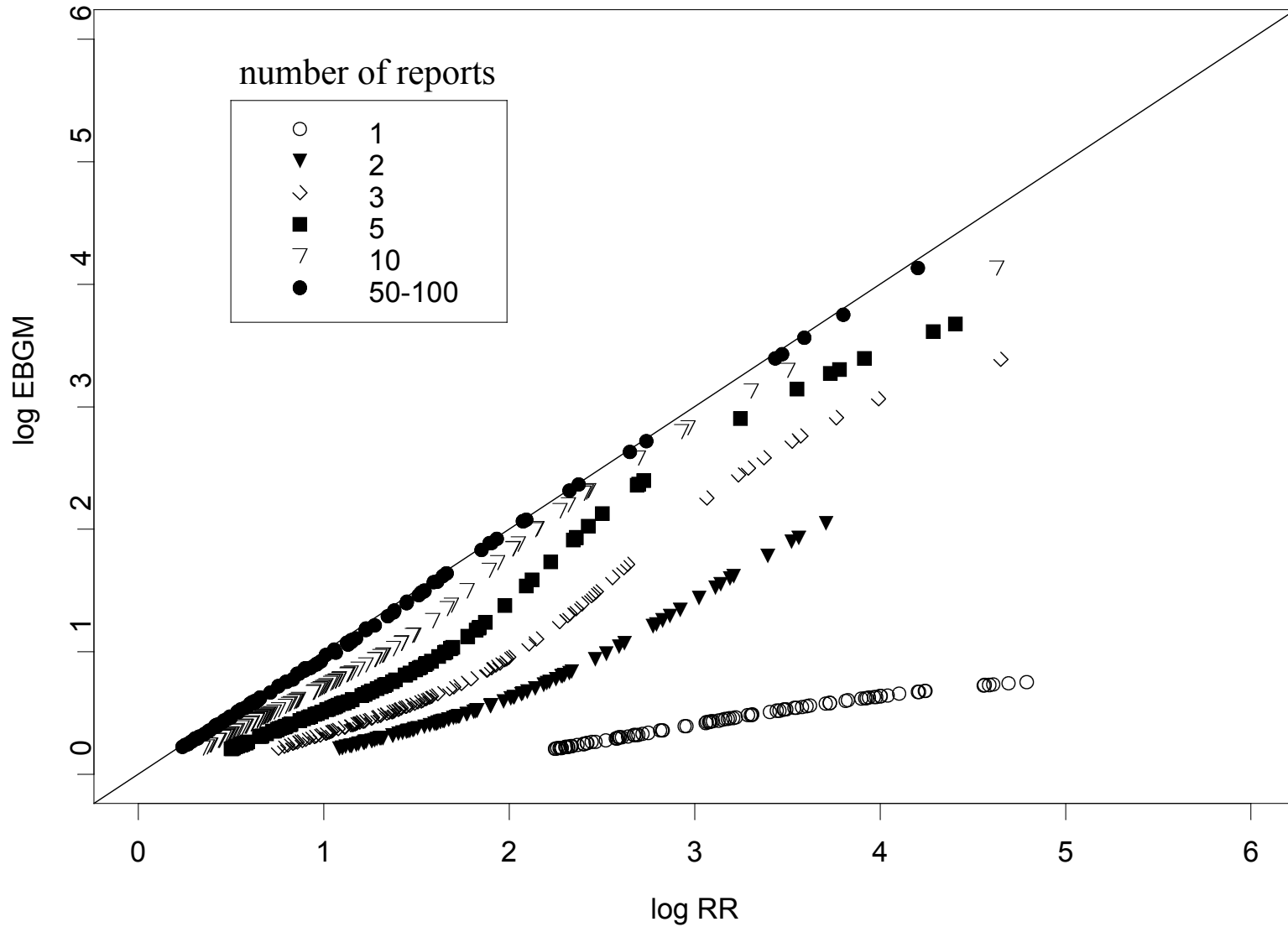
- Disadvantages

- Extreme sampling variability when baseline and observed frequencies are small
($N=1, E=0.01$ v.s. $N=100, E=1$)
- GPS provides a shrinkage estimate of RR that addresses this concern.

$$E_{ij} = N_{ij} * N_{..} / N_{i.} * N_{.j}$$

	AE _j	Not AE _j	
Drug _i	N _{ij}		N _{i.}
Not Drug _i			
	N _{.j}		N _{..}

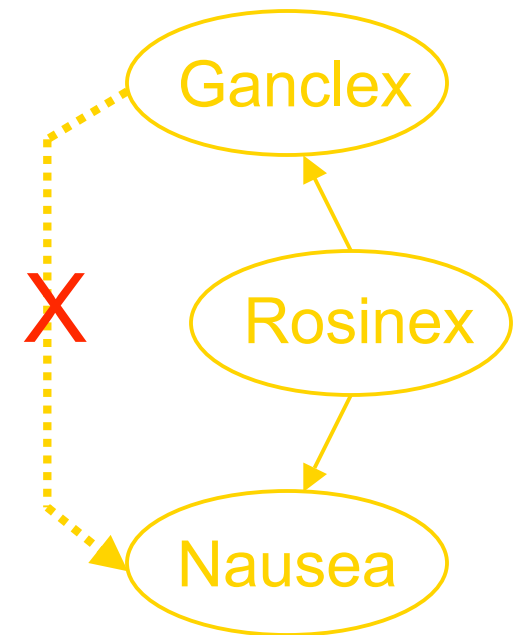
GPS SHRINKAGE – AERS DATA

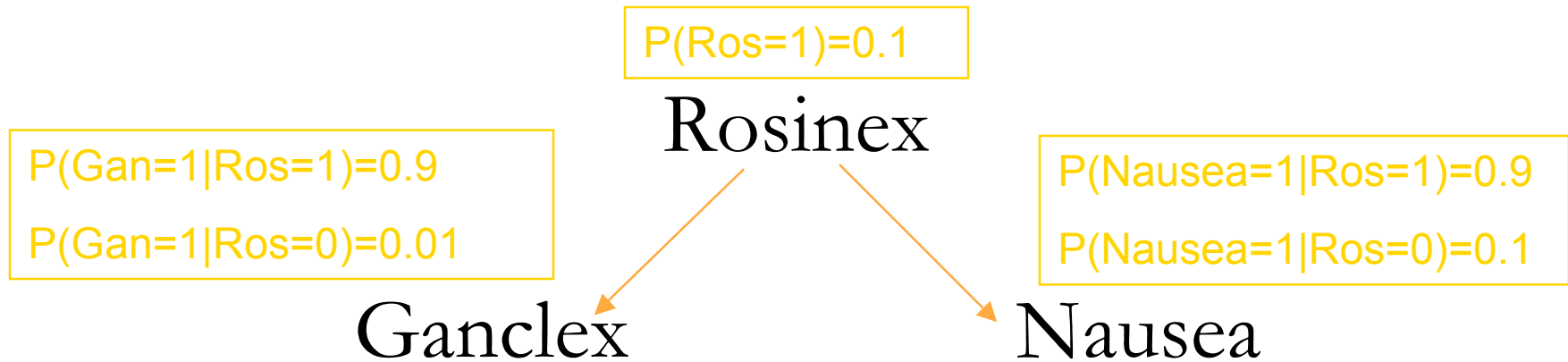


Confounding

- Contingency table analysis ignores effects of drug-drug association on drug-AE association

	Rosinex		No Rosinex		Total	
	Nausea	No Nausea	Nausea	No Nausea	Nausea	No Nausea
Ganclex	81	9	1	9	82	18
No Ganclex	9	1	90	810	99	811
RR	1		1		4.58	





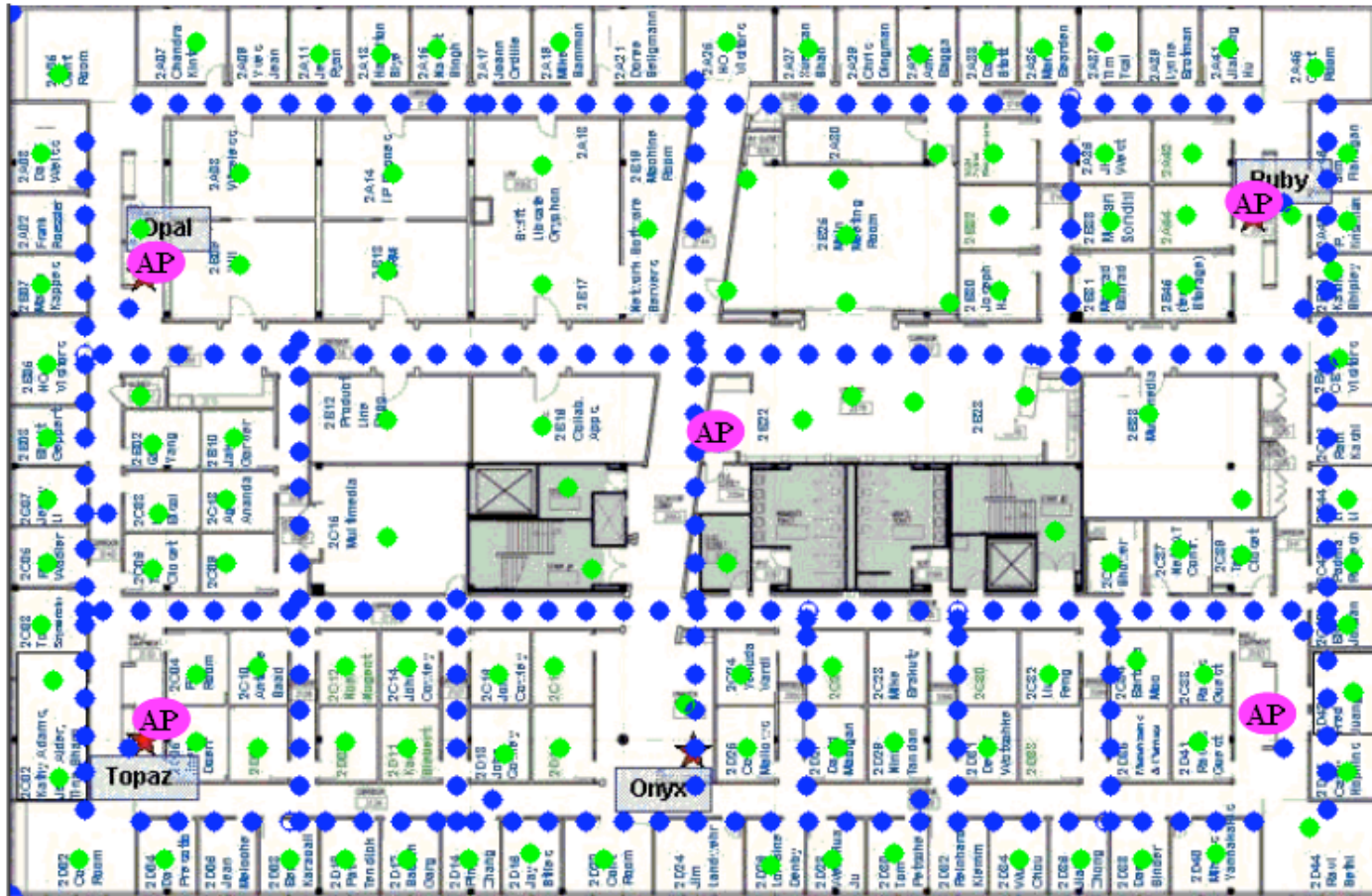
		Nausea vs. Ganclex		Nausea vs. Rosinex	
		Value	Rank	Value	Rank
N		1673	2	1826	1
Bayesian Logistic Regression	Laplace-CV	0.0	9127	4.0	7
GPS EBGM		2.8	73	3.0	68
Observed RR		2.8	744	3.0	681

Logistic Regression

- $\log [P/(1-P)] = \text{intercept} + \sum (\text{each drug effect})$
 - $P = \text{Pr}(\text{report with these drugs will have the AE})$
- 15,000 logistic regressions with $n \approx 3$ million
- 15,000 main effects
- millions of pairwise interactions???

Current Work

- Model associations between *groups* of drugs and *groups* of adverse events
- Bayesian generative approach applicable
- Sketch:
 - assign every drug to a latent group
 - assign every AE to a latent group
 - for each set of drugs and set of AE's, generate a report with probability defined by latent group memberships
- Major computational challenges
- Blei, Feinberg, Ghahramani, Roweis, etc.



- BR has 5 APs, site dimension: 225 ft X 144 ft
 - 259 blue (corridor) data points taken earlier

Hierarchical Model

