

# Education in Statistics

## Time for change? Yes, definitely.

Mattias Villani<sup>1</sup>

**Department of Statistics  
Stockholm University**

Department of Computer and Information Science  
Linköping University



---

<sup>1</sup> Slides: <http://mattiasvillani.com/news>

# My recent article in Qvintensen 2020:1

The cover of the magazine Quintensen features a woman smiling in the background. In the foreground, there is a large, semi-transparent graphic of a human brain with a network of glowing nodes and lines representing connections or data flow. Several small circular icons with symbols like a key, a gear, and a checkmark are connected to different points on the brain's network. The title 'Quintensen' is written in a large, stylized, greenish-blue font at the top. Below it, the issue number 'NR 1 2020' is visible. A small box labeled 'TEMA/' is located in the upper left corner. The main headline 'Artificiell intelligens och maskininlärning' is overlaid on the brain graphic. At the bottom, there is a question in Swedish: 'ÄR DET P-VÄRDEN VI SKA VARA RADDÅ FÖR?'.

# Statistikämnet är hotat – tack och lov

Data finns numera överallt och är en enorm drivkraft för modern industri. I denna blomstringstid för dataanalys finns paradoxt nog en oro för statistikämnets död. Vi utmanas av närliggande ämnesområden som maskinell lämning. Det är en välbefolkt katalysator till förmögelser av värst åttme.

Jag har sedan länge haft  
ett förhållande till natur och teknik  
och dess utveckling. Detta är en  
genetiskt instinkt som jag  
är medveten om. Jag har därför med tidigare  
ärkeförbundit med tekniken och  
med teknologin. Jag har också  
varit med i teknologins utveckling  
och med dess utveckling. Jag  
är medveten om att tekniken  
är en viktig del av den moderna  
samhället.

Detta har inte varit en enskild  
tänkning. Jag har sedan länge haft  
ett förhållande till natur och teknik  
och dess utveckling. Detta är en  
genetiskt instinkt som jag  
är medveten om. Jag har därför med tidigare  
ärkeförbundit med tekniken och  
med teknologin. Jag har också  
varit med i teknologins utveckling  
och med dess utveckling. Jag  
är medveten om att tekniken  
är en viktig del av den moderna  
samhället.

Detta har inte varit en enskild  
tänkning. Jag har sedan länge haft  
ett förhållande till natur och teknik  
och dess utveckling. Detta är en  
genetiskt instinkt som jag  
är medveten om. Jag har därför med tidigare  
ärkeförbundit med tekniken och  
med teknologin. Jag har också  
varit med i teknologins utveckling  
och med dess utveckling. Jag  
är medveten om att tekniken  
är en viktig del av den moderna  
samhället.

ARTIFICIELL INTELLIGENS & MASKINLÆRNING



Sex  
områden  
där jag  
upplever att  
statistik-  
ämnet bör  
förändras

# My background

- PhD in Statistics from Stockholm University, 2000.
- Researcher and adviser Sveriges Riksbank, 2004-11.  
Developed models and methods for monetary policy.
- Interested in Computational Statistics.
- Flexible models and Machine Learning since 2006.
- Built up Division of Statistics and Machine Learning  
within Computer Science dept at LiU, 2011-19.
- Developed a large number of machine learning courses for  
statistics master, engineering programs and industry.<sup>2</sup>
- Application areas in research: economics, neuroscience,  
robotics, software engineering, text analysis, transportation.<sup>3</sup>

---

<sup>2</sup> Slides: <http://mattiasvillani.com/teaching>

<sup>3</sup> Slides: <http://mattiasvillani.com/research>

# 1. Prediction and decision

## ■ Statistics:

- ▶ focus on **estimating** and **interpreting parameters**
- ▶ very little focus on prediction
- ▶ even less on decisions under uncertainty

## ■ Machine learning:

- ▶ focus on **prediction**
- ▶ automated **decisions** are central

*Leo Breiman's (2001) "Statistical Modeling: The Two Cultures." An energetic and passionate argument for the "algorithmic culture" ... Breiman turned out to be more prescient than me: **pure prediction algorithms have seized the statistical limelight in the twenty-first century** ...*

*Brad Efron, "Prediction, Estimation, and Attribution", Discussion paper JASA, June 2020.*

# Let's keep the discipline, but shift our focus

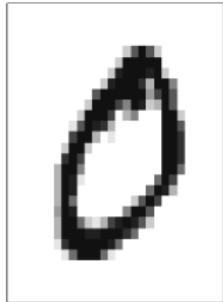
- Not a fan: '**pure prediction algorithms**', e.g. random forests.
- **Flexible probabilistic models** for **prediction** and **decisions**.
- **Attribution** and **Causality** are absolutely crucial  
... but that rarely means  $H_0 : \theta \neq 0$ .
- **Which factors are important for predictions and decisions?**
- **Statistics should insist on** using probability models and well-founded inference methods...
- ... but we need to shift **focus** to stay relevant
- ... especially when **data and computing is everywhere**.

# Prediction and decisions are great for teaching

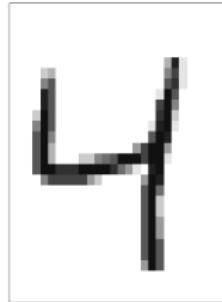
5



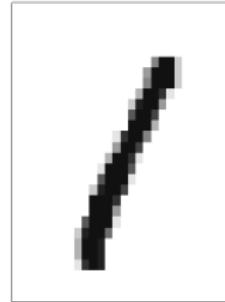
0



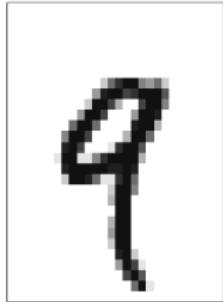
4



1



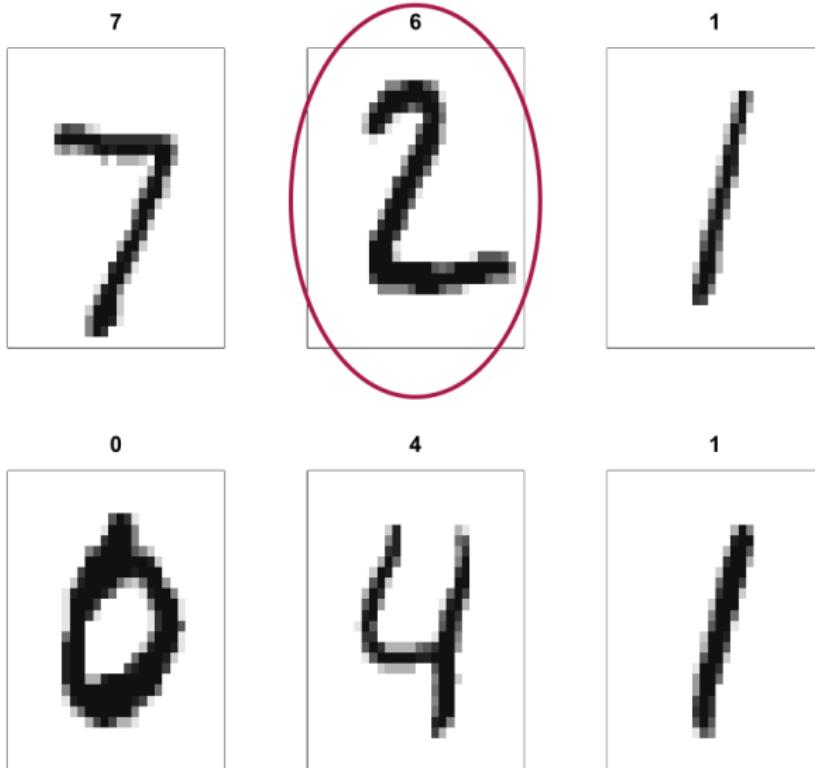
9



2

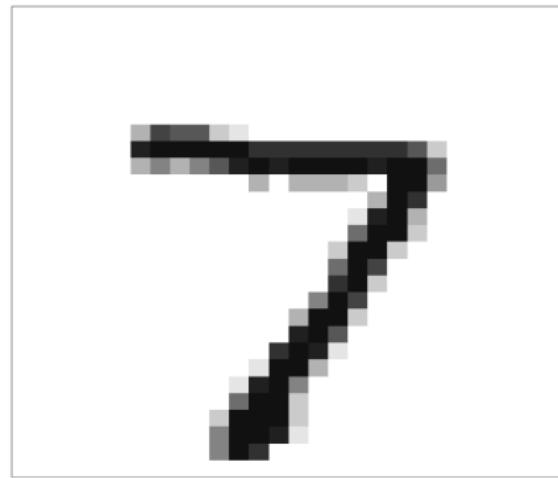


# Predicting new handwritten digits

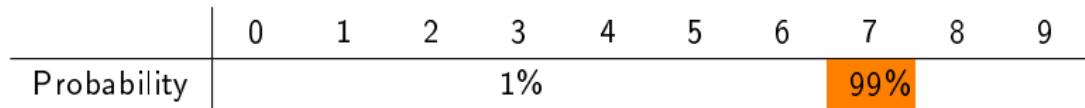


# Probabilities and uncertainty

7

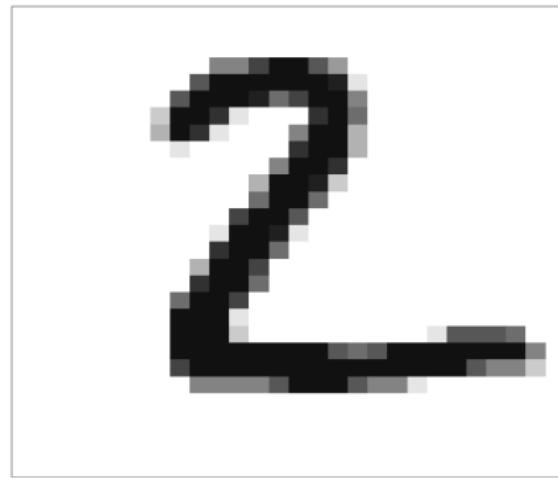


- The system is **very certain** about this prediction

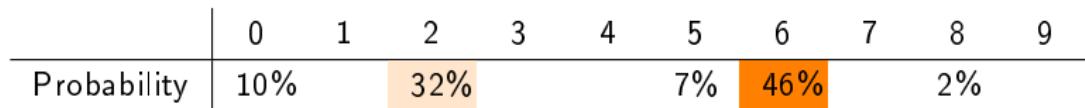


# Probabilities and uncertainty

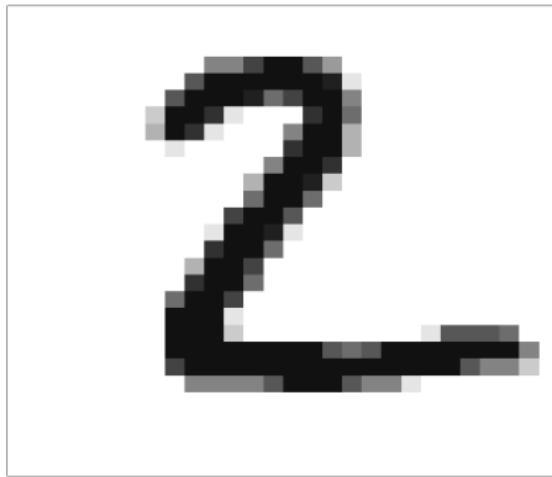
6



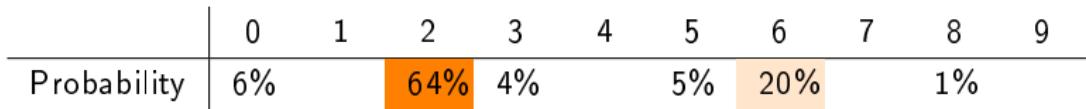
- The system is **very uncertain** about this prediction



# More training data improves performance



- The system is now getting it right



# 1. Prediction and decision

- Above example motivates the study of:
  - ▶ Probability and probability models
  - ▶ Inference
  - ▶ Probabilistic prediction
  - ▶ Decision making
- Complement/replace with example from students' domain.
- Changes in statistics education:
  - ▶ **Much larger focus on prediction.** Early, for motivation.
  - ▶ Compare models' **predictive performance**. Cross-validation.
  - ▶ **Regression early** and a lot of it.
  - ▶ **Less tests.** A lot less.
  - ▶ More focus on **decision making** as the **final aim**.

## 2. Flexible models and regularization

### ■ Flexible models

- ▶ Richly parametrized mean/moments.

Example:  $\mathbb{E}(y|x)$  modelled by spline basis/neural net/tree.

Example:  $\mathbb{V}(y|x)$  is a flexible function of  $x$ .

- ▶ Flexible functional form

Example:  $p(y|x)$  is a mixture density

Example:  $p(y|x)$  is given by a multivariate copula.

- Overfitting is clearly an issue.

- But linear models may underfit.

- **Regularization:** flexible models, but shrink parameters.

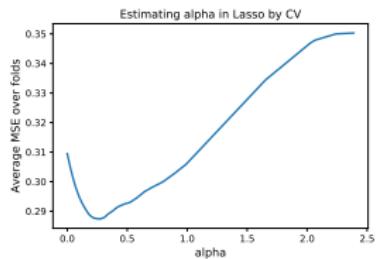
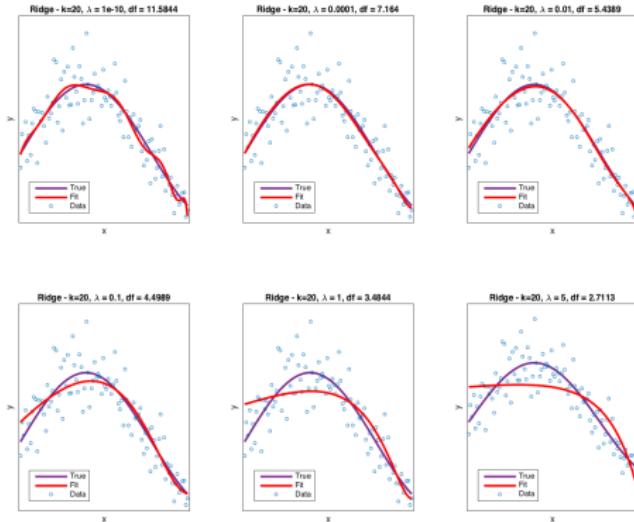
## 2. Flexible models and regularization

- Model:  $y = f_{\theta}(x) + \varepsilon$ ,  $\varepsilon \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$ .
- Maximum likelihood with **L2-regularization** minimizes

$$\sum_{i=1}^n (y_i - f_{\theta}(x_i))^2 + \lambda \|\theta\|^2$$

- Linjär modell:  $f_{\theta}(x_i) = x^{\top} \theta$ , then  $\hat{\theta}_{\text{RR}} = (X^{\top} X + \lambda I_n)^{-1} X^{\top} y$ .
- **L1-regularization: Lasso.** Shrinkage and variable selection.
- Large  $\lambda$ : **Effective flexibility**  $\ll$  nominal number of  $\theta$ 's.
- Linear fit:  $\hat{y} = Hy$ . **Degrees of freedom** of a model:  $\text{tr}(H)$ .
- Problems with  $n < p$  is solved by regularization.

# L2-regularization on polynomial of order 20



## 2. Flexible models and regularization

- Let's teach about:
  - ▶ **Regularization.** Effective flexibility  $\neq$  Nominal flexibility
  - ▶ **Predictive performance.** **Cross-validation.** **Overfitting.**
  - ▶ **Interpreting nonlinear models.**
  - ▶ **Interpretable predictions** and **decisions**.
  - ▶ **Modern variable selection.** Lasso and beyond.

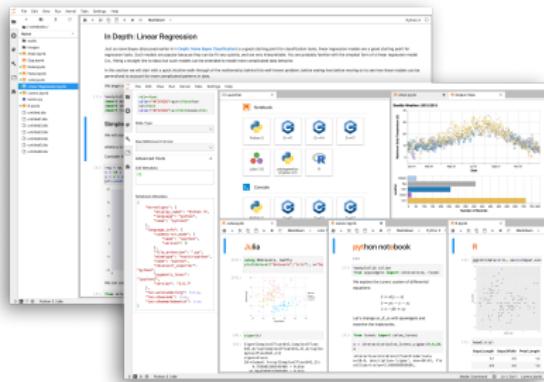
### 3. Computations matter - programming



- My current advice:
  - ▶ learn **R** well
  - ▶ learn **basics of Python** (and/or Julia!)
  - ▶ learn **interoperability** (how to call Python code from R etc)
- Hallmarks of Science:
  - ▶ transparent
  - ▶ open to critique
  - ▶ reproducible
  - ▶ sharing knowledge
- So universities should teach with **open source software**.

### 3. Competing with data scientist - tools matter

- **JuPyteR notebooks:** interactive documents with
  - ▶ **text** (markdown)
  - ▶ **math** (LaTeX)
  - ▶ **code** (>40 languages)
- Great for teaching, exercises and computer exams.



- Package development, Version control (Git), Testing, Profiling.
- Change we must: (**R**) programming early. More **tools**.

## 4. Computational efficiency and large-scale data

- Statistician's **strength**: **statistical efficiency**.
- Statistician's **weakness**: **computational efficiency**.
- **Computational complexity** of a method is often crucial.
- **Large-scale (big) data**
  - ▶ Memory access and storage
  - ▶ Streaming data. Real-time constraints.
  - ▶ **Scalable statistical methods**
- Ch-ch-changes:
  - ▶ Teach awareness that **computations matters**.
  - ▶ **Analysis of algorithms** from computer science.
  - ▶ **Performant code, parallel computing** and all that.

## 5. New data types, noise and anomalies

- Modern data are **not tables**:
  - ▶ Images
  - ▶ Text
  - ▶ Computer generated log files
  - ▶ Sound bites



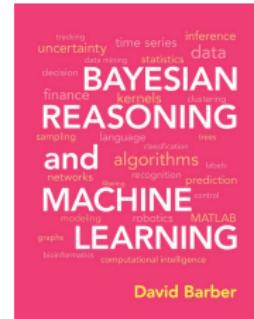
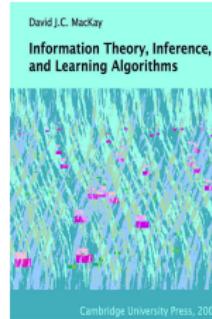
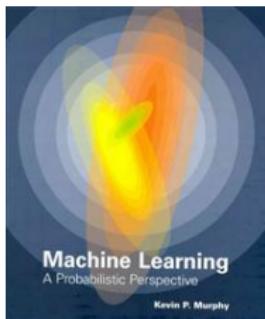
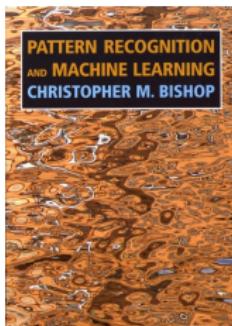
```
w -O3 -fno-passes.o  
code/2_4/compiler_aliases.o  
gcc -w -O3 -fno-strict-aliasing -fcov-  
ode/2_4/compiler_magicsets.o  
gcc -w -O3 -fno-strict-aliasing -fcov-  
ode/2_4/compiler_bitsets.o  
gcc -w -O3 -fno-strict-aliasing -fcov-  
ode/2_4/compiler_importer.o  
code/2_4/compiler_aliases.o  
fno-strict-aliasing -fcov-
```



- Let's get messy:
  - ▶ More **examples** and exercises with **non-standard data** types
  - ▶ **Case studies** with messy, noisy, and slightly inappropriate data
  - ▶ **Data access tools** and **pre-processing** pipelines

## 6. and yes, Bayes is popular in Machine Learning

- Many popular course books in ML take a Bayesian approach.



- Why is Bayes so popular in ML?
  - More **engineering** than pure science. **Subjective** is OK.
  - Using more information (**priors**) is better.
  - Prediction** and **decision-making** are natural for Bayes.
  - Bayes has **regularization** built in (smoothness priors).
  - Algorithmic** implementations (MCMC etc).
- Deep learning - Bayes to **quantify uncertainty**.

# 6. Bayes is popular in applied statistics

## Current issue of JASA Applications & Case Studies:

Journal of the American Statistical Association

Submit an article ■ New content alert ■ RSS ■ Subscribe ■ Citation search

Current issue ■ Browse list of issues ■ Explore

Applications and Case Studies

Article A Bayesian General Linear Modeling Approach to Cortical Surface fMRI Data Analysis >  
Amanda F. Meleg, Yu (Ryan) Yau, David Bolis, Finn Lindgren & Martin A. Lindquist  
Pages 512-523  
Published online: 12 Jun 2019  
Abstract | Full Text | References | PDF (2069 KB) | Supplemental

Article Robust Clustering With Subpopulation-Specific Deviations >  
Diana K. Stepanova, Amy H. Herring & Andrew Orlitzky  
Pages 524-534  
Published online: 18 Jun 2019  
Abstract | Full Text | References | PDF (2688 KB) | Supplemental

Article A Large-Scale Constrained Joint Modeling Approach for Predicting User Activity, Engagement, and Churn With Application to Freemium Mobile Games >  
Trinath Banerjee, Gourab Mukherjee, Shantanu Datta & Pulak Ghosh  
Pages 535-554  
Published online: 18 Jun 2019  
Abstract | Full Text | References | PDF (2677 KB) | Supplemental

Article Hierarchical Space-Time Modeling of Asymptotically Independent Exceedances With an Application to Precipitation Data >  
Jean-Noël Bacro, Carlo Gaetan, Thomas Opitz & Gélydyk Toussenelle  
Pages 555-568  
Published online: 18 Jun 2019  
Abstract | Full Text | References | PDF (2739 KB) | Supplemental

Article Testing and Estimation of Social Network Dependence With Time to Event Data >  
Lin Su, Wentian Lu, Kai Song & Danyang Huang  
Pages 579-592  
Published online: 19 Jun 2019  
Abstract | Full Text | References | PDF (2485 KB) | Supplemental

Article Bayesian Optimal Design for Ordinary Differential Equation Models With Application in Biological Science >  
Anthony M. Oreshnikoff, David C. Woods & Ben M. Parker  
Pages 593-606  
Published online: 25 Jun 2019  
Abstract | Full Text | References | PDF (1464 KB) | Supplemental

Article MIMIX: A Bayesian Mixed-Effects Model for Microbiome Data From Designed Experiments >  
Neal S. Granahan, Yavien Guan, Brian J. Reich, Elizabeth T. Borer & Kevin Gross  
Pages 607-620  
Published online: 20 Jul 2019  
Abstract | Full Text | References | PDF (2308 KB) | Supplemental

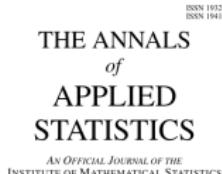
Article Bayesian Graphical Compositional Regression for Microbiome Data >  
Jialiang Mao, Yuhuan Chen & Li Ma  
Pages 621-634  
Published online: 20 Aug 2019  
Abstract | Full Text | References | PDF (2140 KB) | Supplemental

Article Statistical Topology and the Random Interstellar Medium >  
Rutha G. Barnes, Inna Maiselova, Paul Boddy, Andrew Fletcher & Anvar Shukurov  
Pages 635-650  
Published online: 20 Aug 2019  
Abstract | Full Text | References | PDF (1674 KB) | Supplemental

Volume and Issue Information

# 6. Bayes is popular in applied statistics

## ■ Current issue of *Annals of Applied Statistics*:



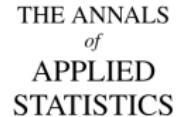
AN OFFICIAL JOURNAL OF THE  
INSTITUTE OF MATHEMATICAL STATISTICS

Articles	
SHOOPER: A probabilistic model of consumer choice with substitutes and complements	FRANCIS P. REIF, SUZANNE ASHES AND DAVID M. BIELE
BART with surgical checkups: An analysis of patient-specific offsets	JORDI E. STERLING, JORDI S. MURRAY, CARLOS M. CARVALHO
Integrating survival analysis with uncertain event times in application to a suicide risk study	YUAN YANG, JUN YAN
Efficient posterior inference of the shape of prior posterior: How to do it?	PAUL J. HERRILL, LOREN WERNICKE, BRIAN D. M. TOLI, LEONARD HELD, CHRISTIAN KUEBLER, CHRISTIAN DUMONT AND DAVID M. D'ANGELLO
Modeling survival distributions and dynamics with bivariate Hawkes processes	BEATRIZ D. MARTIN, DANIELA WITTEN AND AMY D. WILLETT
A statistical analysis of smoking initiation	ABHAY CHAKRABORTY, SOUMENDRA NATH LAHIRI AND ALYSSON WILSON
Surface temperature anomalies in India predicted by a function-valued change-point analysis	YUZHENG ZHENG, XIAO DU, RAO LI AND JIANG ZHENG
Assessing wage state transitions and migration using quantile transition regression	YUZHENG ZHENG, XIAO DU, RAO LI AND JIANG ZHENG
Tfitter: A powerful translation and weighting procedure for combining $\alpha$ -values	HONG ZHANG, TIEGEN TONG, JOHN LANDESK AND ZHEYANG WU
Multidimensional $\chi^2$ -square and the CMB test for non-Gaussianity	KERSTIN SPITZER, MARINA PELEZIĆ AND ANDREA FUTCHER
Mathematical Bayesian models for predicting cognitive impairment from brain evolutionary relationships	MOHAMMAD ELMASRI, MAXWELL J. FARRELL, T. BORISHEM DAVISON AND DAVID A. STEPHENS
Bayesian factor models for probabilistic gene expression profile clustering	TRUONG KUNHAMA, ZERANG RICHARD LI, SAMUEL J. CLARK, JEFFREY L. COOPER, JONATHAN M. LEVINE AND JONATHAN WU
Estimating the health effects of environmental hazards using Bayesian hierarchical regression and sparsity inducing priors	JOSEPH ANTONELLI, MARIANO MARCHETTI, DAVID BRESNAHAN, JONATHAN WU AND JONATHAN LEVINE
Robust Wright and Brant tests	ROBERT WRIGHT AND BRANT COOK
Feature selection for generalized varying coefficient multi-effect models with applications to obesity GWAS	WANGJIAN CHU, YIQUAN LIN, JIANG LIU AND MINGHUA REN
Optimal asset allocation with conjugate Bayesian dynamic linear models	JAROD D. FISHER, DUCROF PETTINGAZZO AND CARLOS M. CARVALHO

Continued on back cover

Vol. 14, No. 1—March 2020

ISSN 1932-6157 (print)  
ISSN 1941-7330 (online)



AN OFFICIAL JOURNAL OF THE  
INSTITUTE OF MATHEMATICAL STATISTICS

Articles—Continued from front cover	
Modeling wild life ignition origins in southern California using latent network prior processes	MATTHEW UPPALA AND MARK S. HANCOCK
Regression for copula-linked compound distributions with applications in modeling aging and mortality	YUZHENG ZHENG, XIAO DU, RAO LI AND JIANG ZHENG
Estimating and forecasting the smoking attributable mortality fraction for both smokers jointly in over 60 countries	YICHENG LI AND ADRIAN E. RATTNER
Measuring human activity spaces from GPS data	YICHENG LI, YUN-CHI CHEN AND ADRIAN DORIA
Curves and surfaces for non-Gaussian data	ZHONGHUA LIU, JIN BARNETT AND XIRONG LIN
A comparison of principal component methods between multiple phenotype regression and multiple SNP regression in genome-wide association studies	ZHONGHUA LIU, JIN BARNETT AND XIRONG LIN
Estimating causal effects in studies of human brain function: New models, methods and criteria	ALICE DUNAI, ELENA MATEICKA, JON GOLDBECK
A hierarchical dependent Dirichlet process prior for modeling bird migration patterns in the UK	ALICE DUNAI, ELENA MATEICKA, JON GOLDBECK
Bayesian mixed effects models for zero-and-one compositions in microbiome data analysis	BOYU ZHANG, JIANG LIU, RICCARDO BACALLINA, STEPHEN P. FALKOWSKI, TOMMI VATTANI, CURTIS HETTRICHOWSKI AND LORENZO TRIPPLI

Continued

Sensitivity analysis for an unobserved moderator in RCT-to-target population generalization of treatment effects

TRUNG QUY KHUONG NGUYEN AND ELIZABETH A. STUART

# Summary

- **Prediction** and **decision making**
- **Regularization**
- **Programming** and computational tools matter
- **Computational complexity** is a statistical issue
- **Modern data** are not nice tables.
- **(Bayes)**

Jo, jag menar det!

