

The paper

frontiers in Applied Mathematics and Statistics

ORIGINAL RESEAR published: 27 November 2

Contain and

Better Autologistic Regression

Mark A. Wolters*

variables, spatial statistics

1. INTRODUCTION

Shanghar Center for Mathematical Sciences, Fusian University, Shanghar, China

Autologistic regression is an important probability model for dichotomous random variables observed along with covariate information. It has been used in various fields for analyzing binary data possessing spatial or network structure. The model can be viewed as an extension of the autologistic model (also known as the Ising model, quadratic exponential binary distribution, or Boltzmann machine) to include covariates. It can also be viewed as an extension of logistic regression to handle responses that are not independent. Not all authors use exactly the same form of the autolooistic regression model. Variations of the model differ in two respects. First, the variable coding-the two numbers used to represent the two possible states of the variables-might differ Common coding choices are (zero, one) and (minus one, plus one). Second, the model might appear in either of two algebraic forms: a standard form, or a recently proposed centered form. Little attention has been paid to the effect of these differences, and the literature shows ambiguity about their importance. It is shown here that changes to either coding or centering in fact produce distinct, non-nested probability models Theoretical results, numerical studies, and analysis of an ecological data set all show that the differences among the models can be large and practically significant. Understanding the nature of the differences and making appropriate modeling choices can lead to significantly improved autologistic regression analyses. The results strongly suggest that the standard model with plus/minus coding, which we call the symmetric autologistic model, is the most natural choice among the autologistic variants.

Keywords: probabilistic graphical models, Markov random fields, logistic regression, correlated binary random

OPEN ACCESS

Edited by: Hou-tieng Wu, Duke University, United States

Reviewed by: John Hughes, John Hughes, United States Antonio Calcagni, University of Dents, Baly

> "Correspondence: Mark A. Wolleys

mwohirs@hudan.edu

Specially section

Mathematics of Computation and Easter Science, a section of the Journal Promities in Applied Mathematics and Statistics

Received: 22 March 2017 Accepted: 10 November 2017 Published: 27 November 2017

Citation Waters MA (2017) Better Autologistic Pegression. Pront. Appl. Math. Stat. 3:24. The autologies (A) model is a probabilistic graphical model for multivariate bioary data. It was introduced to the statical literature by lenge (1); and has also been developed by Kaira et al. (2) and frame the break or the state of the s

When binary responses are observed along with covariate information, the autologistic model may be extended to become the autologistic regression (ALR) model. This model can be viewed as a natural extension of ordinary logistic regression to handle cases where responses are not independent. Under the ALR model, the remomes follow an autologistic distribution, and the

Frontiers in Applied Mathematics and Statistics | www.trontiersin.org

November 2017 J Volume 3 Lacticle 2

2017. Better Autologistic Regression, Frontiers in Applied Mathematics and Statistics

What is "data science?"

JOURNAL OF COMPUTATIONAL AND GRAPHICAL STATISTICS 2017, VOL. 20, NO. 4, 745–766 https://doi.org/10.0080/0018600.2017.1384734 Taylor & Francis

ARTICLE HISTOP

KEYWORDS Cross-study analysis; Data analysis; Data science; Meta

(A) Campus-wide initiatives at NYU, Columbia, MIT, ...
 (B) New master's degree programs in data science, for exam

There are new announcements of such initiatives weekly.²

Many of my audience at the Tukey Centennial-where these

remarks were originally presented-are applied statisticians, and

consider their professional career one long series of exercises in

the above "...collection, manarement, processing, analysis, visu-

alization, and interpretation of vast amounts of heterogeneous data associated with a diverse array of ... applications." In fact,

some presentations at the Tukey Centennial were exemplary narratives of "...collection, management, processing, analysis,

versity of Illinois,

2. Data Science "Versus" Statistics

ple, at Berkeley, NYU, Stanford, Carnegie Mellon, Uni-

Received August 205 Revised August 207

analysis; Predictive modeling; Quantitative programming environments

OPEN ACCESS

50 Years of Data Science

David Donoho

Department of Statistics, Stanford University, Standford, CA

ABSTRACT

More than 50 years ago, John Tukey called for a reformation of academic statistics. In "The Future of Data Analysis," he pointed to the existence of an as-yet unrecognized science, whose subject of interest was Interpret for potential or falls analysis' for to 20 years ago, labo strate models, leff Wu, BII Cleveland, and learning from data, or falls analysis' for to 20 years ago, labo strates, leff Wu, BII Cleveland, and Leo Breman independently once again urged academic statistics to expand its boundaries beyond the classical domain of theoretical statistics; chambers called for more emphasis on data preparation and the states of the states and the states of the states and the states of th presentation rather than statistical modeling: and Breiman called for emphasis on realiction rather than Inference. Cleveland and Wu even suggested the catchy name 'data science' for the envisioned field. A recent and growing phenomenon has been the environment of the science' programs and implicit university, including UC Betweey, NTU, MT, and most prominently, the University of Michigan, which in Soptember 2053 announced a \$2053 mounced as \$2050 mounced as a science histiative' that arms to heir \$5 more faculty. Texting in these new science is that are science as a science in the science in the science in the science in the science is the science in the science in the science is the science is the science in the science in the science is the programs has significant overlap in curricular subject matter with traditional statistics courses wet many academic statisticians perceive the new programs as "cultural appropriation." This article reviews some incred ents of the current 'data science moment,'including recent commentary about data science in the popular media, and about how/whether data science is really different from statistics. The now-contemplated field of data science amounts to a superset of the fields of statistics and machine learning, which adds some technology for "scaling up" to "big data." This chosen superset is motivated by commercial rather than intellectual developments. Choosing in this way is likely to miss out on the really important intellectual event of the next 50 years. Because all of science itself will soon become data that can be mined, the imminent revolution in data science is not about mere "scaling up," but instead the emergence of scientific studies of data analysis science-wide. In the future, we will be able to predict how a proposal to change data analysis workflows would impact the validity of data analysis across all of science, even predicting the impacts field statistics and machine learning than today's data science initiatives, while being able to accommodate the same short-term goals. Based on a presentation at the Tukey Centennial Workshop. Princeton, NJ. September M.

1. Today's Data Science Moment

In September 2015, as I was preparing these remarks, the University of Michigan announced a \$100 million "Data Science Initiative" (DSI)¹, ultimately hiring 35 new faculty.

The university's press release contains bold pronouncements 'Data science has become a fourth approach to scientific discovery,

Into science has become a journ approach to sciency, aucorery, in addition to experimentation, modeling, and computation," said Prorent Marthu Pollack

The website for DSI gives us an idea what data science is:

"This coupling of scientific discovery and practice involves the collection, nonnegement, processing, analysis, visualitation, and interpretation of user anounts of heterogeneous data succidated with a diverse array of scientific, translational, and inter-disciplinary applications."

This announcement is not taking place in a vacuum. A number of DSI-like initiatives started recently, including

CONTACT David Donoho C donohojistanford.edv

For a compension or aboreviations used in this article, see Table 1.
For an updated interactive geographic map of degree programs, see http://data-science-university-programs.silk.co

 You an updated interactive geographic map or degree programs, see might used science unit a 2020/0140 control

to an output locking the second process of t

Donoho (2017, JCGS), 50 Years of Data Science:

- "The value of technical work is judged by the extent to which it benefits the data analyst, either directly or indirectly." [Quoting Cleveland (2001)]
- "...a litmus test re Statistical theorists: do they 'care about the data analyst' or do they not?"

What is "data science?"

Donoho's 6 components of "greater data Science"

The activities of GDS are classified into six divisions:

- 1. Data Gathering, Preparation, and Exploration
- 2. Data Representation and Transformation
- 3. Computing with Data
- 4. Data Modeling
- 5. Data Visualization and Presentation
- 6. Science about Data Science

The paper: autologistic regression (1/13)

Example: *H. vulgaris* data (Carl & Kühn, 2007, *Ecological Modeling*; Bardos et al. 2015 arXiv)



z: presence/absence



 $\Pr(Z_i = 1 | x_i),$ logistic regression

 \mathbf{x}_1 : altitude

The paper: autologistic regression (2/13)

- Dichotomous data with predictors
- Local/spatial association
- The applications involve data on a grid (more generally, graph)

The *autologistic regression (ALR) model* is a pairwise *Markov random field* (MRF) of dichotomous random variables, with a linear predictor.

Applications: ecology, computer vision, dentistry, anthropology, materials science, \ldots

The paper: autologistic regression (3/13)

Let ${\bf Z}$ be a vector of dichotomous random variables.

- Autologistic (AL) model is an MRF.
- Undirected graph with adjacency matrix **A**.
- PMF:

$$f_{\mathbf{Z}}(\mathbf{z}) \propto \exp\left(\mathbf{z}^T \boldsymbol{\alpha} + \frac{1}{2} \mathbf{z}^T \boldsymbol{\Lambda} \mathbf{z}\right)$$

unary term pairwise term

• Let $\pi_i = \Pr(Z_i = \text{high} \mid \text{neighbours})$. Conditional form of the model:

$$\log\left(\frac{\pi_i}{1-\pi_i}\right) = \alpha_i + \sum_{j\sim i} \lambda_{ij} z_j$$

- Let $\alpha = \mathbf{X}\boldsymbol{\beta}$ \Rightarrow autologistic regression
- Let $\Lambda = \lambda A$ \Rightarrow "simple" form of the model

UBC, Aug 3, 2018



The paper: autologistic regression (4/13)

1. The STANDARD model: $\mathbf{Z} \in \{0, 1\}^n$



The paper: autologistic regression (5/13)

2. The **CENTERED** model ($\mathbf{Z} \in \{0, 1\}^n$)

• Traditional model has a problem

- Fix β , increase λ , you will find Z = 1 everywhere.
- Why? Because $\sum_{i \sim i} z_j$ is never negative.
- Caragea & Kaiser (2009): "centered parametrization":

$$\log\left(\frac{\pi_i}{1-\pi_i}\right) = \mathbf{x}_i^T \boldsymbol{\beta} + \sum_{j \sim i} \lambda_{ij} (z_j - \boldsymbol{\mu}_j), \quad \text{where} \quad \boldsymbol{\mu}_j = \frac{e^{\mathbf{x}_j^T \boldsymbol{\beta}}}{1 - e^{\mathbf{x}_j^T \boldsymbol{\beta}}}$$

• μ_j is the independence expectation of the Z_j

The paper: autologistic regression (6/13)

3. The **SYMMETRIC** model

- The responses are *categorical*. Don't have to use $\{0, 1\}$ coding.
- In general, could use $\{\ell, h\}$.
- If ${\bf Z}$ has support $\{\ell,h\}^n,$

$$\mathbf{Y} = a\mathbf{Z} + b\mathbf{1},$$
 where $a = \frac{H-L}{h-\ell}, \quad b = L - a\ell$

has support $\{L, H\}^n$.

The symmetric model is the standard model, with $\mathbf{Z} \in \{-h, h\}^n$

- No centering
- $-\,$ Coding symmetric around 0 $\,$
- We shouldn't change coding without thinking...

The paper: autologistic regression (7/13)

Say $\mathbf{Z} \in \{\ell, h\}^n$, with $f_{\mathbf{Z}}(\mathbf{z}) \propto g(\mathbf{z}; \boldsymbol{\theta})$, but we want our model to use coding $\{L, H\}$.

$$\frac{\text{The right way}}{\mathbf{Y} = a\mathbf{Z} + b\mathbf{1}} \iff \mathbf{Z} = \frac{1}{a}\mathbf{Y} - \frac{b}{a}\mathbf{1}$$
$$f_{\mathbf{Y}}(\mathbf{y}) = \Pr(\mathbf{Y} = \mathbf{y})$$
$$= \Pr(a\mathbf{Z} + b\mathbf{1} = \mathbf{y})$$
$$= f_{\mathbf{Z}}(\frac{1}{a}\mathbf{y} - \frac{b}{a}\mathbf{1})$$
$$\propto g(\frac{1}{a}\mathbf{y} - \frac{b}{a}\mathbf{1}; \boldsymbol{\theta})$$
$$\propto g(\mathbf{z}; \boldsymbol{\theta}).$$

The tempting way

Just plug in $\mathbf{y} = a\mathbf{z} + b\mathbf{1}$. Let the parameter be $\boldsymbol{\theta}'$.

$$\begin{array}{rcl}
f'_{\mathbf{Y}} & \propto & g(\mathbf{y}; \boldsymbol{\theta}') \\
 & \propto & g(a\mathbf{z} + b\mathbf{1}; \boldsymbol{\theta}')
\end{array}$$

To achieve $f'_{\mathbf{Y}} = f_{\mathbf{Y}}$, we need θ' to compensate for linear transformation of \mathbf{z} .

The paper: autologistic regression (8/13)

• Derive the model for arbitrary $\{\ell, h\}$ coding, we find

logit (Pr(
$$Z_i = h | \mathbf{Z}_{-i})$$
) = $(h - \ell) \left[\mathbf{x}_i^T \boldsymbol{\beta} + \sum_{j \sim i} \lambda_{ij} (z_j - \mu_j) \right]$

where

 $\mu_j = \begin{cases} \mathbf{0} & \text{for a standard model} \\ \\ \frac{\ell e^{\ell \alpha_i} + h e^{h \alpha_i}}{e^{\ell \alpha_i} + e^{h \alpha_i}} & \text{for a centered model} \end{cases}$

• Negpotential function:

$$Q(\mathbf{z}) = \mathbf{z}^T \mathbf{X} \boldsymbol{\beta} - \mathbf{z}^T \boldsymbol{\Lambda} \boldsymbol{\mu} + \frac{1}{2} \mathbf{z}^T \boldsymbol{\Lambda} \mathbf{z}$$

• With this, we can study the effect of coding & centering changes.

The paper: autologistic regression (9/13)

Refer to any particular choice of coding and centering as a *variant* of the model. Are all variants equivalent?

Theorem 1: All AL variants are equivalent to any standard model.

Theorem 2: ALR variants are not equivalent, in general.

→ Many variants, all called "autologistic regression models," are actually different, non-nested distribution families.

➡ Exception: all symmetric models are equivalent.

(*Equivalence*: parameter settings always exist that make the two models assign the same probabilities to every configuration of \mathbf{Z} .)

The paper: autologistic regression (10/13)

Theorem 3: Only the symmetric models have reasonable largeassociation behaviour.

"Simple" model. Let λ increase.

- Centered variants behave *counterintuitively* when λ large.
- Symmetric variants are the *only ones* with reasonable behaviour as $\lambda \to \infty$.

The paper: autologistic regression (11/13)



The paper: autologistic regression (12/13)

H. vulgaris fitted models' marginal probabilities

traditional

centered



| | ŝ | | | | c. | el. | 'n, |
|------|---|-----|---|---|----|-----|-----|
| - 20 | | | | | | | |
| - 21 | | | | | | | |
| 5 | | | | | | | |
| 100 | | 100 | - | - | ÷. | | |

symmetric

| Model | $\hat{\beta}_0$ (SE) | $\hat{\beta}_1$ (SE) | $\hat{\lambda} (SE)$ |
|-------------|----------------------|----------------------|----------------------|
| logistic | 2.78 (0.10) | -0.79(0.028) | |
| traditional | -2.12(0.22) | -0.16(0.026) | $1.43 \ (0.066)$ |
| centered | -1.74(0.31) | -0.17(0.040) | $1.51 \ (0.050)$ |
| symmetric | 0.50(0.11) | -0.13 (0.029) | 1.43(0.071) |

The paper: autologistic regression (13/13)



Referee reaction

- Main conclusions: $\{-1, 1\}$ coding is much preferred over $\{0, 1\}$.
 - The centered model has fundamental problems.
 - \rightarrow Most prior ALR analyses are questionable.
 - \rightarrow We need to change the standard of practice.

What were referee reactions?

Referee A "Lovely" "Surprising"... "One of the best papers on the subject."

Referee reaction

- Referee B agreed that: The results do not appear in the literature
 - The results are correct
 - The $\{-1, 1\}$ model is superior



BUT, recommended rejection "both for its unmotivated purpose of study and for its lack of technical sophistication."

After revision:

- It was claimed (six times) that the paper lacked *"intellectual merits"*.
- The use of *"precalculus"* math was listed (four times) as indication of the paper's low quality.
- In two places it was suggested that the paper reflects my lack of understanding of the state of the art

|--|--|--|--|--|--|

Referee reaction

Some quotes:

"There is a difference between creative research and a collection of *mathematically correct but trivial, easy-to-obtain results* that expand into the volume of a research article."

"... this paper is a *summarization of trivial facts supported by unsophisticated mathematics* intentionally expanded to create the illusion of undeserved mathematical complication."

"In the reviewer's career, it is rare to witness a paper with the quality of the current manuscipt published on professional academic journals. ... the reviewer considers the outcomes of publishing re-explanations of common knowledge in shallow mathematics—even without objective error due to the simplistic nature of the technical arguments—catastrophic, since it will prevent original and innovative research from reaching their target readers."

So how did it get published?

Fortunately, *Frontiers in Applied Mathematics and Statistics* has progressive policies.



- Open access.
- Rapid process.
- Review is structured, focused on correctness.
- Can communicate directly with reviewers.
- Reviewers' identity published with the paper.

Between two worlds?

"Big data"/"data science"/"analytics":

- CS/EE culture
- proceedings
- \bullet prediction
- software; utility

"traditional statistics":

- Math culture
- traditional journals
- estimation & inference
- theory, generality, rigor

Possible remedies?

My *personal* plan: 1. Give up trying to impress the math culture.

- 2. Instead, focus on actual impact, within my abilities.
- 3. Hope that academia's attitudes about performance catch up.

How to change focus to "actual impact"?

- Broaden publication targets. All papers are random access anyway.
 - Open access
 - Rapid, correctness-focused review
 - Newer, data-science focused journals
- Alternative outputs
 - Quality software! (all my papers: 51 citations; my 2 R packages: 200+ downloads per month)
 - Publish data sets
- Collaborative, applied work... solve actual problems.

The big challenge

Implementation is hard. Software development takes time.

The basic formula for academic reputation:

- Where you work
- How many papers
- In which journals

Moore-Sloan Data Science Environments (NYU, UCB, UW... http://msdse.org/) "dramatically advance data-intensive scientific discovery..."

"data science in research universities requires precisely the kind of complex, long-term interdisciplinary work with methodological and engineering efforts that leads to low performance under traditional metrics and slow progress and lack of fit in existing career tracks."