## Autologistic Regression Models, With Application to Segmentation of Hyperspectral Satellite Imagery

Mark Wolters, Shanghai Center for Mathematical Sciences Charmaine Dean, Western University

> November 27, 2015 Shanghai Jiao Tong University Institute of Natural Sciences

## Outline

#### The technical content:



Autologistic regression for binary variables C, with complex association, and covariate information X.

- There are a few model variants.
- Claim: *it matters* which one you choose.

#### The application:



Where is the smoke in this picture?

# **Motivating application**

## Remote sensing for smoke monitoring

Earth-orbiting satellites help study large-scale environmental phenomena.

Our interest: smoke from forest fires.





Blue

Original image

Green

T

True class

#### Data: MODIS images

- 1 per day, 143 days
- 1.2 Mp each
- Centered at Kelowna, BC
- Hand-drawn smoke areas

**Goal:** classify pixels into *smoke/nonsmoke* 

#### Why?

- Health studies
- Model input or validation
- Monitoring & archiving

## **Data characteristics**

- *Binary responses* (smoke/nonsmoke).
- Spectra at each pixel are *covariates* for predicting smoke.
- Hyperspectral images: a high-dimensional predictor space.
- Expect *spatial association*.



## Notation



# Autologistic regression models

## **Spatial Associations**

 $Image \ segmentation = pixel \ classification.$ 

If independent pixels  $\Rightarrow$  Use standard classification technology.

But smoke/nonsmoke regions are *spatially smooth*.



#### RGB working image

#### The true scene



## **Model-Based Approaches**

Many ad hoc ways to let pixels influence each other.

Model-based approach: Markov random fields (MRFs).

- graphical model
- popular in computer vision



We will use the discriminative approach.

Model  $p(\mathbf{C}|\mathbf{X}, \theta)$  directly as a MRF.

#### Hammersley-Clifford theorem

Joint PMF can be expressed as a product of *potential functions*, one for each *maximal clique*.

 $\mathcal{M}=$  the set of maximal cliques.  $\mathbf{C}_m=$  the variables in the  $m^{\mathrm{th}}$  clique. Then

$$p(\mathbf{c}) = \frac{1}{Z} \prod_{m \in \mathcal{M}} \phi_m(\mathbf{c}_m)$$

Customary to write the joint density as *Gibbs distribution* form,

$$p(\mathbf{c}) \propto e^{Q(\mathbf{c})}$$

where  $Q(\cdot)$  is the *negpotential function*.



- A Markov random field of **binary random variables**.
- For now, use zero/one coding:  $\mathbf{C}_i \in \{0, 1\}.$
- The graph is a **regular**, **square grid**.
- Nothing new Physics: Ising model Statistics: Besag (1974)

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 1 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 1 0 0 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0 0 1 0 0 0 1 1 1 0 0 0 0 0 0 0 0 0 01100000000000 . . . . . . . . . . . . 0 00000000100000 00000000011110000 0 0 0 0 0 0 0 0 0 1 1 1 1 0 0 0 0 0 0 0 0 0 0 1 1 1 1 0 0 0 0 00100000001100000 <u>0 1 0 0 0 0 0 0 0 0 1 0 0 0 0 0</u> 0 0 0 0 0 0 0 0 0 0 0 0 1 1 0 0 0 0 1 1 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 0 0 0 0 0 0 0 0 0 0 0

#### Joint PMF:

$$\Pr(\mathbf{C} = \mathbf{c} | \boldsymbol{\alpha}, \boldsymbol{\lambda}) = \frac{1}{Z(\boldsymbol{\alpha}, \boldsymbol{\lambda})} \exp\left(\sum_{i \in \mathcal{V}} \alpha_i c_i + \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} c_i c_j\right)$$
  
**unary** terms  
(one per vertex) **pairwise** terms  
(one per edge)

- Positive  $\alpha_i$  values favor +1 class.
- Setting  $\lambda_{ij} > 0$  favors *locally smooth* configurations ( $C_i = C_j$ ).
- Typically set  $\lambda_{ij} = \lambda, \forall i, j$

#### **Conditional distributions:**

Let  $\pi_i = \Pr(C_i = 1 | \text{all other } C)$ . Then can show:

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

$$j \text{ is a neighbour}$$
of  $i$ 

## Autologistic Regression

Put covariates in the unary part:  $\alpha_i = \mathbf{x}_i^T \boldsymbol{\beta}$ .

Pairwise coefficients:  $\lambda_{ij} = \lambda$ .

Then:

$$\log\left(\frac{\pi_i}{1-\pi_i}\right) = \mathbf{x}_i^T \boldsymbol{\beta} + \lambda \sum_{j \sim i} c_j.$$

Interpretation:

- Unary part is a linear predictor.
- $\mathbf{x}_i^T \boldsymbol{\beta}$  determines conditional log-odds of  $C_i = +1$  in the absence of spatial effects.
- Pairwise  $\lambda$  determines strength of neighbour effects.
- Setting  $\lambda = 0$  reverts to standard logistic regression.

## A centered version

Caragea & Kaiser (2009)

"Autologistic models with interpretable parameters"

Hughes, Haran, & Caragea (2011)

"Autologistic models for binary data on a lattice"

- The neighbour sum  $\sum c_j$  increases log-odds unless all neighbours are zero.  $\implies$  Estimates of  $\beta, \lambda$  are strongly coupled
- Strongly recommended a *centered autologistic model*:

$$\log\left(\frac{\pi_i}{1-\pi_i}\right) = \mathbf{x}_i^T \boldsymbol{\beta} + \lambda \sum_{j \sim i} (c_j - \mu_j)$$

where  $\mu_j = E[Y_j | \lambda = 0]$  is the independence expectation.

Estimation is made hard by the normalizing constant.

#### **Existing possibilities**

- 1. Ignore spatial association (logistic regression, large n, large p).
- 2. Pseudolikelihood (PL):  $L(\beta, \lambda) \approx \prod_{i = 1}^{n} \operatorname{logit}(\pi_i)$
- 3. Monte Carlo ML
- 4. Bayesian approach

Hughes et al. use perfect sampling; recommend PL for large n.

#### Problems

- We have  $\sim 10^8 \ {\rm pixels}$
- We have thousands of predictors, need model selection
- We're still developing models—rapid evaluation of candidates is beneficial

# Some claims about coding and centering

## **Model variants**

- Binary variables do not have to take values 0 and 1.
- In general, let them have coding  $\{\ell,h\}$
- We're most interested in

 $\{0,1\}$ , used in statistics and sometimes in computer vision  $\{-1,1\}$ , used in physics and sometimes in computer vision

• And we have two possibilities: centered, or standard (not centered)

Are all these model variants just parameter transformations, or are they distinct models?

If not the same, what are the differences?

In the following, say  $f_1(\mathbf{z}; \boldsymbol{\theta}_1)$  and  $f_2(\mathbf{y}; \boldsymbol{\theta}_2)$  are *equivalent* if for any  $\boldsymbol{\theta}_2$  there exists a  $\boldsymbol{\theta}_1^*$  such that the two models assign the same probability whenever  $\mathbf{z}$  and  $\mathbf{y}$  represent the same configuration.

If we let the coding be  $\{\ell,h\}$  and define the centering adjustment

Then the *autologistic (AL)* PMF is

$$f_{\mathbf{C}}(\mathbf{c}; \boldsymbol{\alpha}, \lambda) \propto \exp(\mathbf{c}^T \boldsymbol{\alpha} - \lambda \mathbf{c}^T \mathbf{A} \boldsymbol{\mu}_{\boldsymbol{\alpha}} + \frac{\lambda}{2} \mathbf{c}^T \mathbf{A} \mathbf{c}),$$

where  ${\bf A}$  is the adjacency matrix.

And the logit form is

$$\log\left(\frac{\pi_i}{1-\pi_i}\right) = (h-\ell)\left(\alpha_i + \lambda \sum_{j\sim i} (c_j - \mu_j^{\alpha})\right)$$

To obtain the *autologistic regression (ALR)* model, just plug in  $\alpha = \mathbf{X} \boldsymbol{\beta}$ 

#### **Theorem 1**

All AL models are equivalent, irrespective of coding or centering.

#### **Sketch of proof:**

• If z and y have different coding, they are related by

 $\mathbf{z} = a\mathbf{y} + b\mathbf{1}$ 

whenever they represent the same configuration.

- $\bullet\,$  Use this to equate PMFs of z and y
- Obtain an explicit parameter transformation from one model to the other

#### **Theorem 2**

Different ALR models are not, in general, equivalent.

#### **Sketch of proof:**

- Follow logic of Theorem 1.
- Transformation between model only exists if an overdetermined system (n equations, p unknowns) can be solved for  $\beta$ .
- Coefficients of the system depend on arbitrary X.
- System is linear for standard models, nonlinear for centered ones.

This means

- standard, zero/one
- standard, plus/minus
- centered, zero/one
- centered, plus/minus

Are four different probabilistic structures.

## Advantages of standard, plus/minus model

**Claim:** the standard model with  $\{-1, 1\}$  coding is the best choice. **Why?** 

1. It resolves the asymmetry of the standard model without the awkward  $\mu_{\alpha}$  term:

$$\log\left(\frac{\pi_i}{1-\pi_i}\right) = 2\left(\mathbf{x}_i^T\boldsymbol{\beta} + \lambda \sum_{j\sim i} c_j\right)$$

Sign of pairwise term depends on majority vote.

- 2. It "decouples"  $\beta$  and  $\lambda$  better than centering (evidence to follow).
- 2. It allows a convenient *plug-in estimation* of  $\lambda$  ...

#### **Proposal: plug-in estimation**

- a) Use independence (logistic) to get  $\hat{oldsymbol{eta}}$ 
  - Including model selection
  - Sample pixels if neccesary to reduce  $\boldsymbol{n}$  to manageable size
- b) Choose  $\hat{\lambda}$  to optimize predictive power

#### Rationale

Treat  $\lambda$  as a smoothing parameter.

- Assuming independence,  $\hat{\beta}$  captures how information in X can be used to predict C.
- For fixed  $\hat{\beta}$ , tuning  $\lambda$  will optimally reduce noise in the predicted probabilities.

# **Checking the claims**

## A small example

n = 9 variables, square graph.



One predictor plus intercept.

- Small problem.
- Can compute probabilities directly.

Linear predictor vector is

$$\boldsymbol{\alpha} = \mathbf{X}\boldsymbol{\beta} = \begin{bmatrix} 1 & x_o \\ 1 & x_5 \\ 1 & x_o \end{bmatrix}$$

- Choose values of  $x_o, x_5, \beta_0, \beta_1$
- Get marginal Pr(C<sub>5</sub> = high) as a function of λ, for four models:
  - standard, zero/one
  - standard, plus/minus
  - centered, zero/one
  - centered, plus/minus

#### Small example: case 1





Shows the asymmetry of the standard 0/1 model.









$$oldsymbol{eta} = \left[ egin{array}{c} 0 \\ 1 \end{array} 
ight], x_o = 0, x_5 = -1.$$
 Linear predictors:  $egin{array}{ccc} 0 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 0 \end{array}$ 



Standard plus/minus very different from the other three.

Limiting probability not equal to 1 or 0.

 $C_5$  probability constant wrt  $\lambda$ .

## **Simulated images**

Generated RGB images with characteristics similar to the smoke data.



- 5 sizes: 100<sup>2</sup>, 200<sup>2</sup>, 400<sup>2</sup>, 600<sup>2</sup> 800<sup>2</sup> pixels.
- 90 images at each size.
- training, validation, and test groups.

Parameter estimates and prediction error: plug-in vs. pseudolikelihood (plus-minus coding).

pixels	method	Ŕ	$\hat{G}$	$\hat{B}$	λ	error rate (%)
100 <sup>2</sup>	plug-in	-2.21	-2.02	1.91	0.90	20.1
	PL	-2.04	-1.99	2.06	0.99	20.4
200 <sup>2</sup>	plug-in	-1.64	-1.35	1.71	1.00	17.7
	PL	-1.61	-1.30	1.70	1.19	17.7
400 <sup>2</sup>	plug-in	-2.05	-1.42	1.63	1.60	20.1
	PL	-2.08	-1.40	1.68	1.36	20.1
600 <sup>2</sup>	plug-in	-1.91	-1.22	1.76	1.95	20.6
	PL	-1.97	-1.36	1.79	1.51	20.4
800 <sup>2</sup>	plug-in	-1.55	-1.44	1.58	1.95	18.8
	PL	-1.57	-1.43	1.49	1.59	18.6

- Parameter estimates similar.
- Error rates similar.

## Simulated images-results 2

Example predictions,  $800 \times 800$  image.





## **Simulated images-results 3**

What about plugin with different model variants?



## Analysis of the smoke data

## **Analysis flowchart**



## **Results 1**

Use a logistic GAM

- Each variable or interaction is a piecewise linear function
- Model search by genetic algorithm

Predictor set	Selected variables (MODIS band numbers)	plug-in $\hat{\lambda}$
main effects	1 6 7 8 14 16 17 18 21 23 25 26 30 31 32 36	1.85
main effects &	7 30 2:3 5:26 6:11 7:36 8:20 8:22 8:25 8:31	1.75
interactions	13:15 13:23 16:31 18:23 22:36 32:36	

	E	rror rate (%	)
Model	nonsmoke	smoke	11
Widdel	pixels	pixels	overall
main effects, logistic	21.1	25.9	21.6
interactions, logistic	20.0	23.3	20.3
main effects, autologistic	17.6	23.9	18.2
interactions, autologistic	16.2	21.3	16.7

## **Results 2**

Qualitative results

- Reasonable results in many cases
- Mixed smoke + cloud is still a problem
- Data quality issues (mis-labelled training/test data)



RGB image

autologistic prediction



# **Conclusions and future directions**

### Summary

- ALR is an interesting option for binary-response regression problems with complex associations.
- Thus far, different communities appear to have used different codings by default.
  - But this yields different models!
  - Plus/minus coding is best?
- The centered model has been put forth as the "new default" ALR model
  - Our work casts doubt on this choice.
- We've proposed a computationally-feasible analysis scheme for ALR with large sets of hyperspectral images.

## **Further work**

- First priority: finish assessments of model variants and formalize
- ALR extensions
  - Let the pairwise parameter be  $\lambda(\mathbf{x}_i, \mathbf{x}_j)$ : adaptive smoothing.
  - Autobinomial model
- Related models
  - MRF of Beta RVs to model probabilities directly?
- Other applications
  - Ecological data
  - Network data
  - ... suggestions?