Living in a Sparse World

Ernst Wit Johann Bernoulli Institute University of Groningen

 $e.c.wit@rug.nl \\ http://www.math.rug.nl/~ernst$

12 April 2010

Joint work with Luigi Augugliaro (Palermo)

▲ 御 ▶ → ミ ▶

Numquam ponenda est pluralitas sine necessitate

William Occam (1288-1348) proposed a meta-theory of knowledge: "For nothing ought to be posited without necessity."

Can be interpreted as a

- Aesthetic principle: enhances model interpretability through parsimonious representation
- Pragmatic principle: computability.
- ► A priori information principle: represents expectation about nature of solution.
- Prediction principle: bias-variance trade-off
- Bayesian principle?

Exercise: Predict the next two numbers in the sequence

$$-1, 3, 7, 11, \dots$$

You probably thought 15, 19, using the hypothesis H_1 : Add 4 to the previous number.

But why not -19.9, 1043.8, using the hypothesis

 H_2 : if x is current number, then next number is given by

$$-x^3/11 + 9/11x^2 + 23/11.$$

Relative evidence of two hypotheses given data D:

$$\frac{P(H_1|D)}{P(H_2|D)} = \frac{P(H_1)}{P(H_2)} \frac{P(D|H_1)}{P(D|H_2)}$$

where $P(H_1)/P(H_2)$ could include an "aesthetic Occam" preference.

Simple hypothesis uses two parameters:

(i) α =starting value and (ii) β =increment, assumed randomly drawn between -50 and +50, in which case

$$P(D|H_1) = P(\alpha = -1; \beta = 4) = \left(\frac{1}{101}\right)^2 = 1 \times 10^{-4}.$$

Complex hypothesis uses (arguably) 7 parameters: (i) α =starting value and (ii-vii) $\frac{\beta_1}{\beta_2}$, $\frac{\beta_3}{\beta_4}$, $\frac{\beta_5}{\beta_6}$ in the polynomial, again, assumed randomly drawn between -50 and +50, in which case

$$P(D|H_2) = ... = \frac{1}{101} \left(\frac{4}{101}\frac{1}{50}\right) \left(\frac{4}{101}\frac{1}{50}\right) \left(\frac{2}{101}\frac{1}{50}\right) = 2.5 \times 10^{-12}.$$

So, $P(H_1|D)/P(H_2|D) = 4 \times 10^7$ IN FAVOUR OF SPARSITY

It is known that deletion and amplification of certain parts of DNA plays a role in the severity of breast-cancer.

John Bartlett (Royal Infirmary, Glasgow) wants to use deletion and amplification data on **62** breast cancer patients across **59** genes.

His expectation: a few gene deletions and gene amplifications will affect the severity of the breast cancer (measured as NPI).

$$\mathsf{NPI}_i = \sum_{j=1}^{59} x_{ij}\beta_j, +\epsilon_i$$
 (patient $i = 1, \dots, 62$),

subject to sparsity, i.e. many $\beta_j \approx 0$.

□ > < E > < E</p>

High-dimensional inference. As we want lots of small β , we consider the constraint maximization of

 $I(\beta)$ subject to $||\beta||_q \leq c$,

whose dual is equal to the penalized likelihood

$$I_{\lambda}(\beta) = I(\beta) - \lambda ||\beta||_q^q$$

There are several special cases:

- ▶ q = 2: Ridge regression (1958)
 Under normality leads again to a simple quadratic form.
- ▶ q = 1: Lasso regression (Tibshirani, 1996) No closed form solution.

Geometry of the L_1 penalty = Sparsity



rijksuniversiteit groningen

日本・モート・モ

Lasso applied to DNA deletion/amplification data

Breast cancer survival LASSO



Some stopping rule selects 7 out of 59 genes.

Ernst Wit Living in a Sparse World

Example II. gene regulatory networks



"A collection of DNA segments in a cell which "interact" with each other via their RNA or proteins and with other substances in the cell, thereby governing the rates at which genes in the network are transcribed."

Gaussian Graphical Models



 Absence/presence of an edge indicates conditional independence of the variables:

▶ < 몰 ▶ < 몰 ▶</p>

▶ No edge joining Y and Z \iff Y⊥Z|rest

If U=(X,Y,Z) in a graph are Gaussian, then $U\sim {\it N}(\mu,\Sigma).$

whereby $\Theta = \Sigma^{-1}$ represents conditional independence, i.e.

if $\theta_{Y,Z} = 0$, then Y is independent of Z given rest.



Estimating Θ

The Gaussian profile likelihood is given by

$$I(\Sigma) = -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (y - \hat{\mu})^t \Sigma^{-1} (y - \hat{\mu}),$$

and so

$$\begin{split} l(\Theta) &= \frac{1}{2} \log |\Theta| - \frac{1}{2} (y - \hat{\mu})^t \Theta(y - \hat{\mu}) \\ &= \frac{1}{2} \left(\log |\Theta| - \operatorname{Trace}(S\Theta) \right), \end{split}$$

where S the empirical covariance matrix of X. Then

$$\frac{\delta I(\Theta)}{\delta \Theta} = \Theta^{-1} - S,$$

and so $\widehat{\Theta} = S^{-1}$.

(1日) (1日)

Estimating the graph using the L_1 penalty

Use the *lasso* regularized log-likelihood:

$$\max_{\Theta} \left[\log |\Theta| - \mathsf{Trace}(S\Theta) - \lambda ||\Theta||_1 \right],$$

with score equations:

$$\Theta^{-1} - S - \lambda \mathsf{Sign}(\Theta) = 0.$$

NOTE: Compare this with the Lasso problem:

$$\min_{\beta}(y - X\beta)^t(y - X\beta) + \lambda ||\beta||_1$$

with solution

$$X^{t}X\beta - X^{t}y + \lambda \operatorname{Sign}(\beta) = 0.$$

・ 回 ト ・ ヨ ト ・

An example: Cell signalling network for different λ



Let's return to our old problem:

$$\mathbf{y} = \mathbf{X}eta + \boldsymbol{\varepsilon}$$

By location and scale transformations we can always assume that the covariates are standardized with mean 0 and unit length, and that the response variable has mean 0,

$$\sum_{i=1}^{n} y_i = 0, \quad \sum_{i=1}^{n} x_{ij} = 0, \quad \sum_{i=1}^{n} x_{ij}^2 = 1 \quad \text{for } j = 1, 2, \dots, p.$$

Graphical definition of Least Angle Regression



LAR analysis of the diabetes study



(left) estimates of regression coefficients $\hat{\beta}_j$ plotted versus $\sum |\hat{\beta}_j|$; (right) absolute current correlation as function of LAR step; heavy curve shows maximum current correlation.

Relationship L_1 and LAR

Let $\hat{\beta}$ be the Lasso solution for some λ , with

$$\hat{\mu} = X \hat{oldsymbol{eta}}.$$

This means that $\hat{\beta}$ solves

$$X^{t}X\beta - X^{t}y + \lambda \mathsf{Sign}(\beta) = 0.$$

Then it is easy to show that the sign of any nonzero coordinate $\hat{\beta}_j$ must agree with the sign s_j of the current correlation $\hat{c}_j = \mathbf{x}_j^T (\mathbf{y} - \hat{\boldsymbol{\mu}}),$ $\operatorname{sign}(\hat{\beta}_j) = \operatorname{sign}(\hat{c}_j) = s_j.$ (1)

Therefore, a simple modification LAR produces Lasso estimates:

Remove variable from active LAR set a.s.a. (1) is violated.

□ > < E > < E</p>

IV. Generalized linear models

 ${f Y}$ is a random variable with pdf

$$p_{\mathbf{Y}}(\mathbf{y}; \boldsymbol{\theta}, \lambda) = a(\mathbf{y}; \lambda) \exp\{\lambda(\mathbf{y}^{T}\boldsymbol{\theta} - k(\boldsymbol{\theta}))\}, \quad \mathbf{y} \in \mathcal{Y} \subseteq \mathbb{R}^{n}, \quad (2)$$

with respect to a σ -finite measure ν on \mathbb{R}^n , where $a(\cdot)$ and $k(\cdot)$ are specific given functions.

GLMs (McCullagh and Nelder, 1989):

$$E(Y) = \mu$$

$$\eta = \mathbf{x}^T \boldsymbol{\beta}$$

$$g(\mu) = \eta$$

In order to simplify our notation we denote

$$\mu(\beta) = \left(g^{-1}(\mathbf{x}_1^{\mathsf{T}}\beta), g^{-1}(\mathbf{x}_2^{\mathsf{T}}\beta), \dots, g^{-1}(\mathbf{x}_n^{\mathsf{T}}\beta)\right)^{\mathsf{T}}.$$

Example: logistic regression

Observations Y are binary and Bernoulli distributed,

$$p(y) = p^{y}(1-p)^{1-y}$$
$$= \exp\left(y \operatorname{logit}(p) - \log(\frac{1}{1-p})\right)$$

Therefore,

 $\theta = \text{logit}(p).$

and

$$E(Y) = p$$

$$\eta = \mathbf{x}^T \boldsymbol{\beta}$$

$$\eta = \text{logit}(p)$$

Let

$$\mu(\beta) = \left(\mathsf{logit}(\mathsf{x}_1^{\mathsf{T}}\beta), \mathsf{logit}(\mathsf{x}_2^{\mathsf{T}}\beta), \dots, \mathsf{logit}(\mathsf{x}_n^{\mathsf{T}}\beta) \right)^{\mathsf{I}} \cdot \mathsf{primum restrict strains of the state of th$$

Differential manifold

Let \mathcal{M}^n be the mean value parameter space.

```
Amari (1982):
\mathcal{M}^n can be treated as a n-dimensional differentiable manifold.
```

Vos (1991): $\mu(\beta)$ is an embedding with domain $\mathcal B$

$$\widetilde{\mathcal{M}}^{p}=\{\mu\in\mathcal{M}^{n}\,:\,\mu=\mu(oldsymbol{eta}),\,\, ext{with}\,\,oldsymbol{eta}\in\mathcal{B}\}.$$

The vector space

$$\mathcal{T}_{\boldsymbol{\mu}(\boldsymbol{\beta})}\widetilde{\mathcal{M}}^{\boldsymbol{\rho}} = \mathsf{span}\{\partial_{\beta_1}\boldsymbol{\mu}(\boldsymbol{\beta}), \partial_{\beta_2}\boldsymbol{\mu}(\boldsymbol{\beta}), \dots, \partial_{\beta_p}\boldsymbol{\mu}(\boldsymbol{\beta})\}$$

is called *tangent space* of $\widetilde{\mathcal{M}}^p$ at $\mu(\beta)$.

Metric space using the Fisher information

The expected Fisher information matrix defines an inner product, denoted by $\langle , \rangle_{\mu(\beta)}$, on each tangent space:

$$\langle \mathbf{v}_1, \mathbf{v}_2 \rangle_{p(\mu)} = E_{\mu} \left(\sum_{i=1}^n d\mu_{i,1} \partial_i \ell(\mu) \cdot \sum_{j=1}^n d\mu_{2,j} \partial_j \ell(\mu) \right)$$

= $d\mu'_1 l(\mu) d\mu_2.$

A given generalized linear model can be treated as a Riemannian submanifold of $(\mathcal{M}^n, \langle , \rangle_{\mu(\beta)})$. Let *r* be the *tangent residual vector*

$$\mathbf{r}(\boldsymbol{\mu}(\boldsymbol{\beta})) = \sum_{i=1}^{n} (y_i - \mu_i(\boldsymbol{\beta})) \partial_i \ell(\boldsymbol{\mu}(\boldsymbol{\beta})) \in T_{p(\boldsymbol{\mu}(\boldsymbol{\beta}))} \mathcal{S}$$

It can be shown that the MLE of β is defined as

$$r(\mu(\hat{\beta})) \perp T_{p(\mu(\hat{\beta}))} \mathcal{M}.$$

Angles in a GLM

We will define a solution path in parameter space and need some notation:

let
$$\beta(\gamma)$$
 be a double differentiable curve;
let $\ell(\beta(\gamma)) = \ell(\gamma)$ be the log-likelihood function;
 $\partial_{\beta_m}\ell(\gamma) = \partial_m\ell(\gamma)$ be the derivative of $\ell(\gamma)$ with respect to β_m .

The derivative of the likelihood $\partial_m \ell(\beta(\gamma))$ can be written as inner product between $\mathbf{r}(\beta(\gamma))$ and *m*-th base of tangent space of \mathcal{M} at $\mu(\beta(\gamma))$,

$$\partial_m \ell(\beta(\gamma)) = \left\langle \partial_m \mu(\beta(\gamma)); \mathbf{r}(\beta(\gamma)) \right\rangle_{\mu(\beta(\gamma))}.$$
(4)

Using the law of cosines, this is equivalent with

$$\partial_{m}\ell(\beta(\gamma)) = \cos\left(\rho_{m}(\beta(\gamma))\right) \cdot \|\mathbf{r}(\beta(\gamma))\|_{\mu(\beta(\gamma))} \cdot \|\partial_{m}\mu(\beta(\gamma))\|_{\mu(\beta(\gamma))}$$

= $\cos\left(\rho_{m}(\beta(\gamma))\right) \cdot \|\mathbf{r}(\beta(\gamma))\|_{\mu(\beta(\gamma))} \cdot i_{m}^{1/2}(\beta(\gamma))$



rijksuniversiteit groningen



Ernst Wit Living in a Sparse World

Š

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶

rijksuniversiteit groningen

Extension of LAR (1)

Using expression (4) we have the following differential geometric identity $% \left(\frac{1}{2} \right) = 0$

$$\rho_m(\beta(\gamma)) = \arccos \frac{\partial_m \ell(\gamma)}{\|\mathbf{r}(\beta(\gamma))\|_{\boldsymbol{\mu}(\beta(\gamma))} \|\partial_m \boldsymbol{\mu}(\beta(\gamma))\|_{\boldsymbol{\mu}(\beta(\gamma))}}.(6)$$

where

- $\rho_m(\beta(\gamma))$ is the angle between $\mathbf{r}(\beta(\gamma))$
- ► $\partial_m \mu(\beta(\gamma))$, $i_m(\beta(\gamma))$ is the expected Fisher information for $\beta_m(\gamma)$
- ▶ $\|\cdot\|_{\mu(\beta(\gamma))}$ is the norm defined on the tangent space.

(6) shows that the gradient of the log-likelihood function does not generalize the equiangularity condition, since we are not considering the variation related to

$$i_m^{1/2}(\boldsymbol{\beta}(\boldsymbol{\gamma})) = \|\partial_m \boldsymbol{\mu}(\boldsymbol{\beta}(\boldsymbol{\gamma}))\|_{\boldsymbol{\mu}(\boldsymbol{\beta}(\boldsymbol{\gamma}))}.$$

Expression (6) suggests that a genuine generalization of the Least Angle Regression method should be based on the following indentity

$$r_m^{\mu}(\gamma) = i_m^{-1/2}(\beta(\gamma)) \cdot \partial_m \ell(\gamma) = \cos\left(\rho_m(\beta(\gamma))\right) \cdot \|\mathbf{r}(\beta(\gamma))\|_{\mu(\beta(\gamma))}.$$
(7)



Simulation study

We have evaluated the dgLARS method by means of a simulation study.

We have simulated 1000 samples of (100,200,500) observations with 5000 variables using the following generalized linear model:

$$\begin{array}{rcl} Y_i & \sim & \mathsf{Ber}(\pi_i),\\ \mathsf{logit}(\pi_i) & = & \eta_i,\\ \eta_i & = & 1 + 2 \cdot x_{i1} + 3 \cdot x_{i2} + 4 \cdot x_{i3}. \end{array}$$

- ▶ We have used IID standard Gaussian predictors to obtain the design matrix X.
- The size of the true model, i.e. the number of predictors with non-zero coefficients, is 3.



Figure (a) shows the solution path $\beta_{\mathcal{A}}(\gamma)$ obtained with the proposed algorithm after 15 steps. Figure (b) shows the path of $|\mathbf{r}^{u}(\gamma)|$ as function of the number of variables included in the active set. Heavy lines identify the paths related to the variables used to simulate the response variable.

Comparison L_1 and dgLAR

200 data sets from a logistic regression model. Sample size was equal to 50/100 and the number of covariates was equal to 500/1000.

We considered a model with two important groups as follow

$$(X_1, X_2, ..., X_5) \backsim N_5(\mathbf{0}, \Sigma),$$

 $(X_6, X_7, ..., X_{10}) \backsim N_5(\mathbf{0}, \Sigma),$
 $X_i \backsim N(0, 1), \quad i = 11, ..., p,$

with

$$\Sigma = \begin{pmatrix} 1 & \rho & \dots & \rho & \rho \\ \rho & 1 & \dots & \rho & \rho \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \rho & \rho & \dots & 1 & \rho \\ \rho & \rho & \dots & \rho & 1 \end{pmatrix}$$

and $\rho = 0.9$. We chose $\beta = (5, \dots, 5, 0, \dots, 0)$.

					k/n		
	n	р	0.2	0.4	0.6	0.8	1
dgLARS	50	500	5.04	5.31	5.53	5.7	5.83
			(0.24)	(0.25)	(0.24)	(0.23)	(0.24)
	50	1000	4.92	5.23	5.44	5.62	5.80
			(0.23)	(0.24)	(0.23)	(0.23)	(0.23)
	100	500	5.95	6.29	6.47	6.66	6.80
			(0.20)	(0.20)	(0.20)	(0.20)	(0.19)
	100	1000	5.81	6.04	6.21	6.42	6.61
			(0.22)	(0.23)	(0.22)	(0.22)	(0.21)
L ₁ -penalty	50	500	5.04	5.26	5.21	5.2	5.15
			(0.25)	(0.24)	(0.24)	(0.23)	(0.23)
	50	1000	4.89	5.13	5.15	5.17	5.14
			(0.22)	(0.23)	(0.22)	(0.22)	(0.21)
	100	500	5.99	6.02	6.08	6.03	5.91
			(0.20)	(0.21)	(0.20)	(0.20)	(0.20)
	100	1000	5.81	5.92	5.99	5.92	5.90
			(0.22)	(0.23)	(0.22)	(0.22)	(0.22)

The numbers in parentheses are the corresponding coefficient of deviations.

Average number of true variables identified by the dgLARS and L_1 -regularized logistic regression model based on 200 replicates. The average is expressed as function of the ratio between the steps of the algorithms (k) and the sample size (n).



Figure shows the average number of average number of true variables identified by the dgLARS and L_1 -regularized logistic regression model based on 200 replications. The average is expressed as function of the ratio between the steps of the algorithms (k) and the sample size (n).

ksuniversitei

Conclusions:

- Penalized inference and least angle extensions are important methods to deal with high-dimensional feature spaces.
- L₁ penalized graphical models have interesting genomic applications.
- dgLAR is based on a natural generalization of the geometrical theory underlying the original LAR algorithm.

- [1] Allgower E. and Georg K.(1990) *Numerical Continuation Methods*, Springer, Berlin.
- [2] Amari S.-I.(1985) Differential-Geometrical Methods in Statistics (Lecture Notes in Statistics, 28), Springer-Verlag, New-York.
- [3] Candes E. and Tao T.(2007) The Dantzig Selector: Statistical Estimation when *p* is much larger than *n. Ann. Statist.*, **35**(6), 2313-2351.
- [4] Efron B., Hastie T., Johnstone I. and Tibshirani R.(2004) Least Angle Regression. Ann. Statist., 32(2), 407-499.
- [5] Fan J. and Li R.(2001) Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties. J. Amer. Statist. Assoc., 96(456), 1348-1360.
- [6] Fan J. and Lv J.(2008) Sure independence screening for ultrahigh dimensional feature space. J. R. Statist. Soc. B, **70**(5), 849-911.
- [7] Frank I.E. and Friedmann J.H.(1993) A Statistical View of Some Chemometrics Regression Tools. *Technometrics*, 35(2), 109-135.
- [8] Hastie, T., Tibshirani, R. and Friedmann, J.(2001) The Elements of Statistical Learning: Data mining, Inference, and Prediction, Springer.

イロト イヨト イヨト イヨト

References II

- [9] Jørgensen B.(1987) Exponential Dispersion Models (with discussion). J. R. Statist. Soc. B, 49, 127-162.
- [10] McCullagh P. and Nelder J.(1989) Generalized Linear Models, Chapman & Hall.
- [11] Park M.Y. and Hastie T.(2007) L₁-regularization path algorithm for generalized linear models. J. R. Statist. Soc. B, 69(4), 659-677.
- [12] Rao C.R.(1945) On the distance between two populations. Sankhya, 9, 246-248.
- [13] Tibshirani, R.(1996) Regression Shrinkage and Selection via the Lasso. J. R. Statist. Soc. B, 58(1), 267-288.
- [14] Vos P.W.(1991) A Geometric Approach to Detecting Influential Cases. Ann. Statist., 19(3), 1570-1581.
- [15] Witten, D.M. and Tibshrani, R (2009) Covariance-regularized regression and classification for high-dimensional problems. J. R. Statist. Soc. B, 71(3), 615-636.
- [16] Zou, H. and Hastie, T.(2005) Regularization and variable selection via the elastic net. J. R. Statist. Soc. B, 67(2), 301-320.

riiksuniversiteit

(日) (四) (三) (三) (三)