

# Review of modern logistic regression methods

Enes Makalic Daniel F. Schmidt

Centre for MEGA Epidemiology  
The University of Melbourne

23rd Australasian Joint Conference on Artificial Intelligence  
2010

# Outline

- 1 Introduction
  - Problem Description
  - Motivation
- 2 Regularised Logistic Regression
  - Ridge Regression
  - LASSO
  - Elastic Net
  - Non-negative Garrote
- 3 Simulation
  - Simulated Data
  - Real data

# Outline

- 1 Introduction
  - Problem Description
  - Motivation
- 2 Regularised Logistic Regression
  - Ridge Regression
  - LASSO
  - Elastic Net
  - Non-negative Garrote
- 3 Simulation
  - Simulated Data
  - Real data

# Problem Description (1)

- We have a binary classification problem
  - Data  $\mathbf{y} = (y_1, \dots, y_n)'$ ,  $y_i = \{0, 1\}$
  - Matrix of  $p$  covariate vectors  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$ ,  $\mathbf{x}_j \in \mathbb{R}^n$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

- Use a logistic regression model ( $n$  samples,  $p$  predictors)

## Problem Description (2)

- Logistic regression model for explaining data  $\mathbf{y}$

$$p(y_i = 1 | \mathbf{x}_i, \boldsymbol{\beta}) = \frac{1}{1 + \exp(-\mathbf{x}_i' \boldsymbol{\beta})}$$

- $\boldsymbol{\beta} \in \mathbb{R}^p$  is the vector of logistic regression coefficients
- Task:
  - Estimate parameters
  - Select significant regressors

## Problem Description (2)

- Logistic regression model for explaining data  $\mathbf{y}$

$$p(y_i = 1 | \mathbf{x}_i, \boldsymbol{\beta}) = \frac{1}{1 + \exp(-\mathbf{x}_i' \boldsymbol{\beta})}$$

- $\boldsymbol{\beta} \in \mathbb{R}^p$  is the vector of logistic regression coefficients
- Task:
  - Estimate parameters
  - Select significant regressors

## Problem Description (2)

- Logistic regression model for explaining data  $\mathbf{y}$

$$p(y_i = 1 | \mathbf{x}_i, \boldsymbol{\beta}) = \frac{1}{1 + \exp(-\mathbf{x}_i' \boldsymbol{\beta})}$$

- $\boldsymbol{\beta} \in \mathbb{R}^p$  is the vector of logistic regression coefficients
- Task:
  - Estimate parameters
  - Select significant regressors

# Motivation

- Many problems with maximum likelihood and stepwise regression
- Ideally, want a method that
  - consistently selects true predictors
  - automatically shrinks parameters
  - selects important variables
  - can be applied when  $p \gg n$
  - has the Oracle property (asymptotically)



# Outline

- 1 Introduction
  - Problem Description
  - Motivation
- 2 Regularised Logistic Regression
  - Ridge Regression
  - LASSO
  - Elastic Net
  - Non-negative Garrote
- 3 Simulation
  - Simulated Data
  - Real data

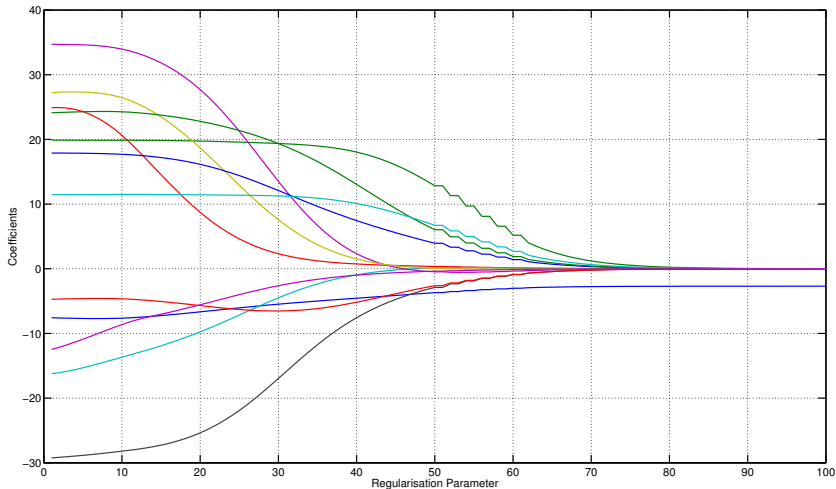
# Ridge Regression (1)

- Ridge regression estimate

$$\hat{\beta}_{\text{RR}} = \arg \max_{\beta} \{l(\beta)\} \quad \text{s.t.} \quad \sum_{j=1}^p \beta_j^2 \leq t$$

- Restricted maximum likelihood estimator
- Introduced originally for linear regression models
  - Dominates least squares in terms of MSE
- Shrinks all parameters to zero ( $t = 0$ ) or includes all predictors ( $t \rightarrow \infty$ )
- Introduces little bias; large reduction in variance

## Ridge Regression (2)



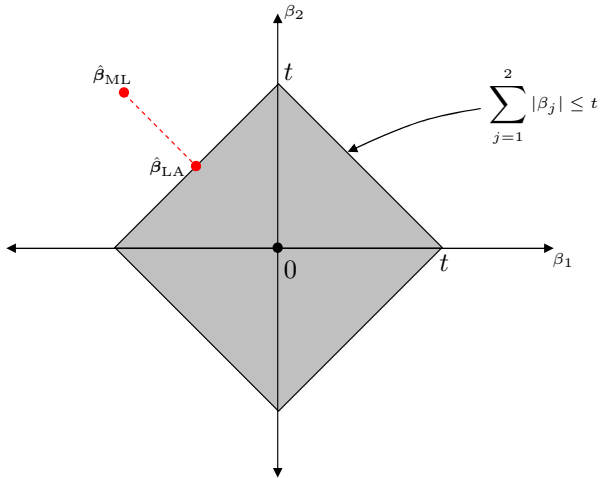
# LASSO (1)

- Least absolute shrinkage and selection operator (LASSO)

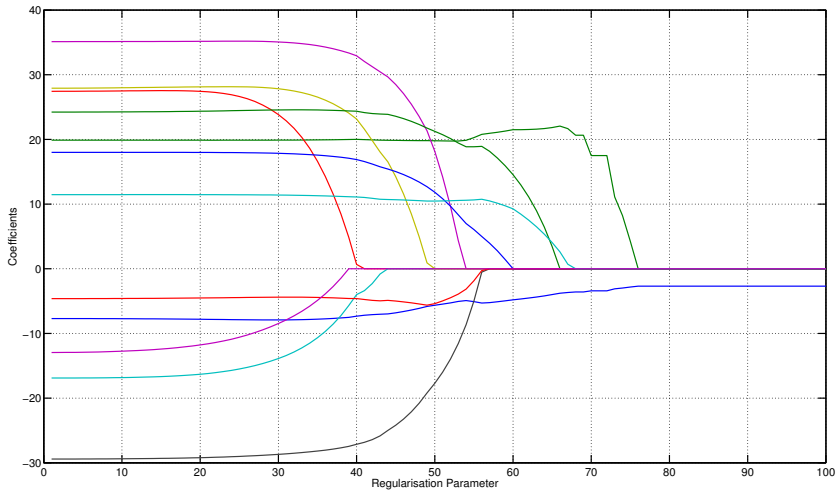
$$\hat{\beta}_{\text{LA}} = \arg \max_{\beta} \{l(\beta)\} \quad \text{s.t.} \quad \sum_{j=1}^p |\beta_j| \leq t$$

- Simultaneous parameter shrinkage and variable selection
- Shrinkage determined by  $t$
- Consistent if  $p \gg n$  under certain assumptions

# LASSO (2)



# LASSO (3)



# Elastic Net

- LASSO can perform poorly if the predictors are correlated
  - Ridge regression performs well
- Solution: combine LASSO and ridge regression penalties

$$\hat{\beta}_{\text{EN}} = \arg \max_{\beta} \{l(\beta)\} \quad \text{s.t.} \quad \sum_{j=1}^p |\beta_j| \leq t_1, \sum_{j=1}^p \beta_j^2 \leq t_2$$

- Handles correlated predictors
- Tends to perform better than LASSO when  $p > n$

# Non-negative Garrote (1)

- Requires an initial parameter estimate  $\beta^*$ 
  - For example, maximum likelihood, ridge regression, etc.
- Non-negative Garrote (NNG) estimate

$$\hat{\beta}_{\text{NG}} = \arg \max_{\tilde{\beta}} \left\{ l(\tilde{\beta}_1, \dots, \tilde{\beta}_p) \right\} \quad \text{s.t. } c_j \geq 0, \sum_{j=1}^p c_j \leq t$$

where  $\tilde{\beta}_j = c_j \beta_j^*$ ,  $j = 1, \dots, p$ .



# Outline

- 1 Introduction
  - Problem Description
  - Motivation
- 2 Regularised Logistic Regression
  - Ridge Regression
  - LASSO
  - Elastic Net
  - Non-negative Garrote
- 3 Simulation
  - Simulated Data
  - Real data

# Simulated Data

- In all simulations ...
  - Sample size  $n = \{20, 50, 100\}$ .
  - Regressor correlation  $\text{corr}(i, j) = 0.5^{|i-j|}$ .
- **Example 1:**  $\beta = (3, 2, 1.5, 0, 0, 0, 0, 0)'$ .
- **Example 2:**  $\beta_j = 0.85$  for all  $j$ .
- **Example 3:**  $\beta = (5, 0.5, 0.5, 0.5, 0, 0, 0, 0)'$ .
- Test metrics
  - median negative log-likelihood (NLL), mean model size (Size), mean number of false positive / false negative regressors

Methods	$n = 20$				$n = 50$			
	NLL	Size	FP	FN	NLL	Size	FP	FN
fwd	30.70 (0.80)	1.24 (0.04)	1.96 (0.03)	0.20 (0.02)	7.12 (0.07)	2.79 (0.04)	0.63 (0.03)	0.42 (0.03)
gfwd	26.64 (0.41)	1.18 (0.04)	1.97 (0.03)	0.15 (0.02)	6.74 (0.05)	2.75 (0.04)	0.64 (0.02)	0.38 (0.03)
lasso	21.13 (0.15)	4.53 (0.05)	0.50 (0.02)	2.03 (0.04)	6.50 (0.04)	5.78 (0.04)	0.05 (0.01)	2.83 (0.04)
glasso	22.40 (0.18)	2.89 (0.04)	0.93 (0.03)	0.82 (0.03)	6.44 (0.05)	3.97 (0.04)	0.23 (0.01)	1.20 (0.04)
rr	21.13 (0.14)	8.00 (0.00)	0.00 (0.00)	5.00 (0.00)	6.76 (0.03)	8.00 (0.00)	0.00 (0.00)	5.00 (0.00)
grr	21.86 (0.21)	3.37 (0.05)	0.77 (0.02)	1.14 (0.03)	6.45 (0.03)	4.32 (0.04)	0.17 (0.01)	1.49 (0.04)
enet	20.64 (0.12)	6.10 (0.05)	0.21 (0.01)	3.31 (0.05)	6.50 (0.03)	6.40 (0.04)	0.02 (0.00)	3.42 (0.04)
genet	22.20 (0.22)	3.07 (0.04)	0.85 (0.02)	0.92 (0.03)	6.42 (0.04)	4.06 (0.04)	0.20 (0.01)	1.26 (0.04)
ilasso	22.41 (0.19)	2.90 (0.04)	0.93 (0.02)	0.83 (0.03)	6.42 (0.05)	3.98 (0.04)	0.22 (0.01)	1.20 (0.04)
nng	21.34 (0.25)	3.50 (0.04)	0.66 (0.02)	1.16 (0.03)	6.34 (0.03)	4.35 (0.04)	0.11 (0.01)	1.46 (0.04)

**Example 1:** median negative log-likelihood (NLL), mean model size (Size), mean number of false positive regressors (FP) and mean number of false negative regressors (FN) included in the selected model. Tests are based on 1000 iterations with standard errors included in parentheses

Methods	$n = 20$				$n = 50$			
	NLL	Size	FP	FN	NLL	Size	FP	FN
fwd	35.10 (0.08)	1.17 (0.05)	6.83 (0.05)	0.00 (0.00)	10.20 (0.09)	4.27 (0.07)	3.73 (0.07)	0.00 (0.00)
gfwd	34.66 (0.03)	1.11 (0.04)	6.89 (0.04)	0.00 (0.00)	9.19 (0.08)	4.05 (0.07)	3.94 (0.07)	0.00 (0.00)
lasso	24.83 (0.19)	4.95 (0.05)	3.06 (0.05)	0.00 (0.00)	7.76 (0.04)	6.94 (0.03)	1.06 (0.03)	0.00 (0.00)
glasso	28.24 (0.19)	3.14 (0.05)	4.86 (0.05)	0.00 (0.00)	8.44 (0.05)	5.66 (0.05)	2.34 (0.05)	0.00 (0.00)
rr	21.23 (0.11)	8.00 (0.00)	0.00 (0.00)	0.00 (0.00)	7.18 (0.03)	8.00 (0.00)	0.00 (0.00)	0.00 (0.00)
grr	26.97 (0.19)	3.80 (0.05)	4.20 (0.05)	0.00 (0.00)	8.23 (0.04)	6.10 (0.04)	1.90 (0.04)	0.00 (0.00)
enet	21.59 (0.12)	7.40 (0.04)	0.60 (0.04)	0.00 (0.00)	7.24 (0.02)	7.88 (0.01)	0.12 (0.01)	0.00 (0.00)
genet	27.20 (0.22)	3.65 (0.05)	4.35 (0.05)	0.00 (0.00)	8.24 (0.04)	6.06 (0.04)	1.94 (0.04)	0.00 (0.00)
ilasso	28.23 (0.21)	3.15 (0.05)	4.85 (0.05)	0.00 (0.00)	8.44 (0.05)	5.67 (0.05)	2.33 (0.05)	0.00 (0.00)
nng	26.49 (0.23)	4.08 (0.05)	3.92 (0.05)	0.00 (0.00)	8.05 (0.04)	6.33 (0.04)	1.67 (0.04)	0.00 (0.00)

**Example 2:** median negative log-likelihood (NLL), mean model size (Size), mean number of false positive regressors (FP) and mean number of false negative regressors (FN) included in the selected model. Tests are based on 1000 iterations with standard errors included in parentheses

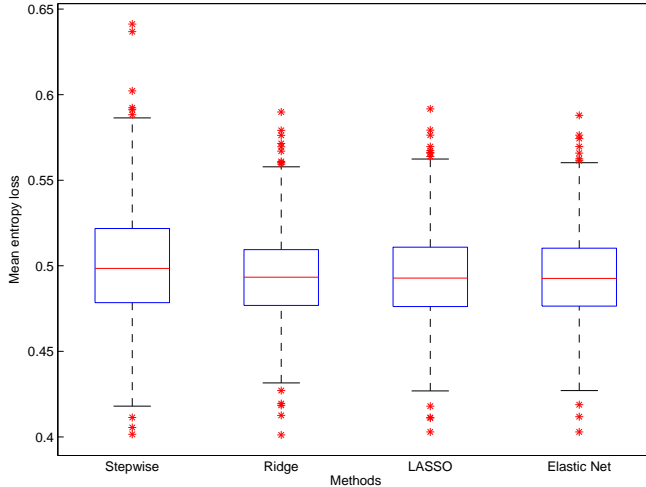
Methods	$n = 20$				$n = 50$			
	NLL	Size	FP	FN	NLL	Size	FP	FN
fwd	19.42 (0.78)	1.17 (0.04)	2.98 (0.03)	0.15 (0.02)	5.57 (0.03)	1.68 (0.04)	2.54 (0.02)	0.22 (0.02)
gfwd	16.55 (0.35)	0.99 (0.03)	3.09 (0.02)	0.08 (0.01)	5.40 (0.02)	1.63 (0.04)	2.57 (0.02)	0.19 (0.02)
lasso	16.70 (0.14)	4.08 (0.05)	1.49 (0.03)	1.57 (0.03)	5.45 (0.03)	5.08 (0.05)	0.98 (0.02)	2.06 (0.04)
glasso	15.63 (0.15)	2.13 (0.03)	2.35 (0.02)	0.48 (0.02)	5.35 (0.02)	2.86 (0.05)	1.90 (0.03)	0.76 (0.03)
rr	20.35 (0.18)	8.00 (0.00)	0.00 (0.00)	4.00 (0.00)	6.02 (0.03)	8.00 (0.00)	0.00 (0.00)	4.00 (0.00)
grr	15.74 (0.11)	2.59 (0.04)	2.12 (0.03)	0.71 (0.02)	5.36 (0.02)	3.21 (0.05)	1.76 (0.03)	0.97 (0.03)
enet	16.70 (0.14)	4.84 (0.06)	1.19 (0.03)	2.03 (0.04)	5.49 (0.03)	5.50 (0.05)	0.83 (0.02)	2.34 (0.04)
genet	15.64 (0.12)	2.21 (0.04)	2.31 (0.02)	0.52 (0.02)	5.34 (0.02)	2.94 (0.05)	1.86 (0.03)	0.80 (0.03)
ilasso	15.63 (0.15)	2.13 (0.04)	2.35 (0.02)	0.48 (0.02)	5.35 (0.02)	2.87 (0.05)	1.89 (0.03)	0.76 (0.03)
nng	15.83 (0.11)	2.76 (0.04)	1.99 (0.03)	0.75 (0.03)	5.30 (0.02)	3.46 (0.05)	1.49 (0.03)	0.95 (0.03)

**Example 3:** median negative log-likelihood (NLL), mean model size (Size), mean number of false positive regressors (FP) and mean number of false negative regressors (FN) included in the selected model. Tests are based on 1000 iterations with standard errors included in parentheses

## Real data

- Obtained from the UCI Machine Learning repository
  - Pima Indian diabetes (250/250/268)
  - Wisconsin diagnostic breast cancer (100/100/369)
- Test metric: mean entropy (classification) loss on test data
- Repeated for 1000 iterations

# Pima indian diabetes



# Wisconsin diagnostic breast cancer

