

Estimation, shrinkage and penalization

Sofia Olhede

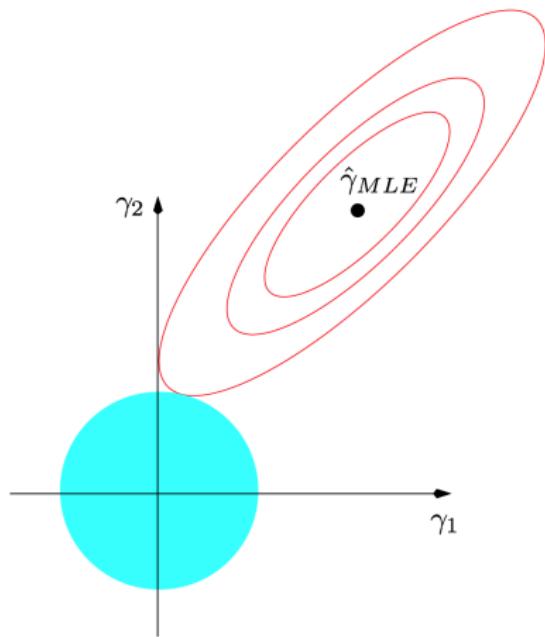


November 16, 2020

1 How and why to shrink?

2 Generalised Linear Model

Shrinkage

Figure: L^2 Shrinkage [Ridge Regression]

Shrinkage II

Theorem

Let $Z_{n \times q}$ be a matrix of rank $r \leq q$ with centred column vectors of unit norm. Given $\lambda > 0$, the unique minimiser of

$$Q(\alpha, \xi) = \|Y - \alpha \mathbf{1} - Z\xi\|_2^2 + \lambda \|\xi\|_2^2$$

is

$$(\hat{\beta}_0, \hat{\gamma}) = (\bar{Y}, (Z^\top Z + \lambda I)^{-1} Z^\top Y).$$

Proof.

Write

$$Y = \underbrace{(Y - \bar{Y}\mathbf{1})}_{=Y^* \in M^\perp(1)} + \underbrace{\bar{Y}\mathbf{1}}_{\in M(1)}$$

Note also that by assumption $\mathbf{1} \in M^\perp(Z)$. Therefore by Pythagoras' theorem

$$\|Y - \hat{\beta}_0 \mathbf{1} - Z\hat{\gamma}\|_2^2 = \left\| \underbrace{(\bar{Y} - \hat{\beta}_0)\mathbf{1}}_{\in M(1)} + \underbrace{(Y^* - Z\hat{\gamma})}_{\in M(Z)} \right\|_2^2 = \|(\bar{Y} - \hat{\beta}_0)\mathbf{1}\|_2^2 + \|(Y^* - Z\hat{\gamma})\|_2^2.$$

Shrinkage III

$$\text{Therefore, } \min_{\alpha, \xi} Q(\alpha, \xi) = \min_{\alpha} \|(\bar{Y} - \alpha)\mathbf{1}\|_2^2 + \min_{\xi} \left\{ \|(\mathbf{Y}^* - \mathbf{Z}\xi)\|_2^2 + \lambda \|\xi\|_2^2 \right\}$$

Clearly, $\arg \min_{\alpha} \|(\bar{Y} - \alpha)\mathbf{1}\|_2^2 = \bar{Y}$ while the second component can be written

$$\min_{\xi \in \mathbb{R}^q} \left\| \begin{pmatrix} \mathbf{Z} \\ \sqrt{\lambda} \mathbf{I}_{q \times q} \end{pmatrix} \xi - \begin{pmatrix} \mathbf{Y}^* \\ \mathbf{0}_{q \times 1} \end{pmatrix} \right\|_2^2$$

using block notation. This is the usual least squares problem with solution

$$\left[(\mathbf{Z}^\top, \sqrt{\lambda} \mathbf{I}_{q \times q}) \begin{pmatrix} \mathbf{Z} \\ \sqrt{\lambda} \mathbf{I}_{q \times q} \end{pmatrix} \right]^{-1} (\mathbf{Z}^\top, \sqrt{\lambda} \mathbf{I}_{q \times q}) \begin{pmatrix} \mathbf{Y}^* \\ \mathbf{0}_{q \times 1} \end{pmatrix} = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^\top \mathbf{Y}^*$$

Note that $\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}$ is indeed invertible. Writing $\mathbf{Z}^\top \mathbf{Z} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top$, we have

$$\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top + \mathbf{U} (\lambda \mathbf{I}_{q \times q}) \mathbf{U}^\top = \mathbf{U} (\boldsymbol{\Lambda} + \lambda \mathbf{I}_{q \times q}) \mathbf{U}^\top$$

and $\boldsymbol{\Lambda} = \text{diag} \underbrace{\{\lambda_1, \dots, \lambda_r\}}_{>0} \underbrace{\{\lambda_{r+1}, \dots, \lambda_q\}}_{=0}$ ($\mathbf{Z}^\top \mathbf{Z} \succeq 0$ & $\text{rank}(\mathbf{Z}^\top \mathbf{Z}) = \text{rank}(\mathbf{Z})$).

To complete the proof, observe that $\mathbf{Z}^\top \mathbf{Y}^* = \mathbf{Z}^\top \mathbf{Y} - \bar{Y} \mathbf{Z}^\top \mathbf{1} = \mathbf{Z}^\top \mathbf{Y}$. □

Shrinkage IV

Note that if the SVD of Z is $Z = V\Lambda U^\top$, last steps of previous proof may be used to show that

$$\hat{\gamma} = \sum_{j=1}^q \frac{\lambda_j}{\lambda_j^2 + \lambda} (V_j^\top Y) U_j,$$

where the V_j s and U_j s are the columns of V and U , respectively.

Compare this to the ordinary least squares solution, when $\lambda = 0$:

$$\hat{\gamma} = \sum_{j=1}^q \frac{1}{\lambda_j} (V_j^\top Y) U_j,$$

which is not even defined if Z is of reduced rank.

Role of λ is to reduce the size of $1/\lambda_j$ when λ_j becomes very small.

Shrinkage V

Proposition

Let $\hat{\gamma}$ be the ridge regression estimator of γ . Then

$$\text{bias}(\hat{\gamma}, \gamma) = -\lambda \left(\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_{q \times q} \right)^{-1} \gamma$$

and

$$\text{cov}(\hat{\gamma}) = \sigma^2 (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1}.$$

Proof.

Since $\mathbb{E}(\hat{\gamma}) = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^\top \mathbb{E}(\mathbf{Y}) = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^\top \mathbf{Z} \gamma$, the bias is

$$\begin{aligned} \text{bias}(\hat{\gamma}, \gamma) &= \mathbb{E}(\hat{\gamma}) - \gamma = \{(\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^\top \mathbf{Z} - \mathbf{I}\} \gamma \\ &= \left\{ \left(\frac{1}{\lambda} \mathbf{Z}^\top \mathbf{Z} + \mathbf{I} \right)^{-1} \left(\frac{1}{\lambda} \mathbf{Z}^\top \mathbf{Z} + \mathbf{I} - \mathbf{I} \right) - \mathbf{I} \right\} \gamma \\ &= \left\{ \mathbf{I} - \left(\frac{1}{\lambda} \mathbf{Z}^\top \mathbf{Z} + \mathbf{I} \right)^{-1} - \mathbf{I} \right\} \gamma = - \left(\frac{1}{\lambda} \mathbf{Z}^\top \mathbf{Z} + \mathbf{I} \right)^{-1} \gamma. \end{aligned}$$

The covariance term is straightforward. □

Shrinkage VI

Corollary (Domination over Least Squares)

Assume that $\text{rank}(Z_{n \times q}) = q$ and let

$$\tilde{\gamma} = (Z^\top Z)^{-1} Z^\top Y \quad \& \quad \hat{\gamma}_\lambda = (Z^\top Z + \lambda I)^{-1} Z^\top Y$$

be the least squares estimator and ridge estimator, respectively. Then,

$$\mathbb{E}\{(\tilde{\gamma} - \gamma)(\tilde{\gamma} - \gamma)^\top\} - \mathbb{E}\{(\hat{\gamma}_\lambda - \gamma)(\hat{\gamma}_\lambda - \gamma)^\top\} \succeq 0$$

for all $\lambda \in (0, 2\sigma^2/\|\gamma\|^2)$.

Ridge estimator uniformly better than least squares! How can this be?
 (What happened to Gauss-Markov?)

- Gauss-Markov only covers unbiased estimators – but ridge estimator biased.
- A bit of bias can improve the MSE by reducing variance!
- Also, there is a catch! The “right” range for λ depends on unknowns.
- Choosing a good λ is all about balancing bias and variance.

Shrinkage VII

Proof.

From our bias/variance calculations on the ridge estimator, we have

$$\begin{aligned} & \mathbb{E}\{(\tilde{\gamma} - \gamma)(\tilde{\gamma} - \gamma)^\top\} - \mathbb{E}\{(\hat{\gamma}_\lambda - \gamma)(\hat{\gamma}_\lambda - \gamma)^\top\} = \\ &= \sigma^2(\mathbf{Z}^\top \mathbf{Z})^{-1} - \sigma^2(\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} + \lambda \left(\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_{q \times q} \right)^{-1} \gamma \\ &= \lambda(\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} \left(\sigma^2(2\mathbf{I} + \lambda(\mathbf{Z}^\top \mathbf{Z})^{-1}) - \lambda \gamma \gamma^\top \right) (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1}. \end{aligned}$$

This last term can be made positive definite if

$$2\sigma^2 \mathbf{I} + \sigma^2 \lambda (\mathbf{Z}^\top \mathbf{Z})^{-1} - \lambda \gamma \gamma^\top \succeq 0.$$

Noting that we can always write

$$\mathbf{I} = \frac{\gamma \gamma^\top}{\|\gamma\|^2} + \sum_{j=1}^{q-1} \boldsymbol{\theta}_j \boldsymbol{\theta}_j^\top$$

for $\{\gamma/\|\gamma\|, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{q-1}\}$ an orthonormal basis of \mathbb{R}^q we see that $\lambda \in (0, 2\sigma^2/\|\gamma\|^2)$ suffices for positive definiteness to hold true.

□

Shrinkage VIII

Role of λ : Regulates Bias–Variance tradeoff

- $\lambda \uparrow$ decreases variance (collinearity) but increases bias
- $\lambda \downarrow$ decreases bias but variance inflated if collinearity exists

Recall:

$$\mathbb{E}\|\hat{\gamma} - \gamma\|^2 = \underbrace{\|\mathbb{E}\hat{\gamma} - \gamma\|^2}_{\text{bias}^2} + \underbrace{\mathbb{E}\|\hat{\gamma} - \mathbb{E}\hat{\gamma}\|^2}_{\text{variance} = \text{trace}[\text{cov}(\hat{\gamma})]} + \underbrace{2(\mathbb{E}\hat{\gamma} - \gamma)^T \mathbb{E}[\hat{\gamma} - \mathbb{E}\hat{\gamma}]}_{=0}$$

Writing $Z^T Z = U \Lambda U^T$ $\text{trace}\{\text{cov}(\hat{\gamma})\} = \sum_{j=1}^q \frac{\lambda_j}{\lambda_j + \lambda} \sigma_j^2$

So choose λ so as to optimally balance **bias/variance**

Use cross validation!



Shrinkage IX

Motivated from Ridge Regression formulation can consider:

$$\begin{aligned} \min! \quad & \| \mathbf{Y} - \beta_0 \mathbf{1} - \mathbf{Z} \boldsymbol{\gamma} \|_2^2 \quad \text{subject to} \quad \sum_{j=1}^{p-1} |\gamma_j| = \|\boldsymbol{\gamma}\|_1 \leq r(\lambda) \\ & \iff \\ \min! \quad & \| \mathbf{Y} - \beta_0 \mathbf{1} - \mathbf{Z} \boldsymbol{\gamma} \|_2^2 + \lambda \|\boldsymbol{\gamma}\|_1. \end{aligned}$$

Shrinks coefficient *size* by different version of *magnitude*.

- Resulting estimator non-linear in \mathbf{Y}
- No explicit form available (unless $\mathbf{Z}^\top \mathbf{Z} = \mathbf{I}$), needs quadratic programming algorithm

Why choose a different type of norm?

L^1 penalty (almost) produces a “continuous” model selection!

Shrinkage X

When the explanatory variables are orthogonal (i.e. $Z^\top Z = I$), then the LASSO⁷ exactly performs model selection via thresholding:

Theorem (Orthogonal Design \leftrightarrow Model Selection)

Consider the linear model

$$\underset{n \times 1}{Y} = \underset{1 \times 1}{\beta_0} \underset{n \times 1}{1} + \underset{n \times (p-1)}{Z} \underset{(p-1) \times 1}{\gamma} + \underset{n \times 1}{\varepsilon}$$

where $Z^\top 1 = 0$ and $Z^\top Z = I$. Let $\hat{\gamma}$ be the least squares estimator of γ ,

$$\hat{\gamma} = (Z^\top Z)^{-1} Z^\top Y = Z^\top Y.$$

Then, the unique solution to the LASSO problem

$$\min_{\beta_0 \in \mathbb{R}, \gamma \in \mathbb{R}^{p-1}} \{ \|Y - \beta_0 1 - Z\gamma\|_2^2 + \lambda \|\gamma\|_1 \}$$

is given by $(\hat{\beta}_0, \check{\gamma}) = (\beta_0, \check{\gamma}_1, \dots, \check{\gamma}_{p-1})$, defined as

$$\hat{\beta}_0 = \bar{Y} \quad \& \quad \check{\gamma}_i = \text{sgn}(\hat{\gamma}_i) \left(|\hat{\gamma}_i| - \frac{\lambda}{2} \right)_+, \quad i = 1, \dots, p-1.$$

Shrinkage XI

Proof (*).

Note that since $Z^\top \mathbf{1} = 0$ and since β_0 does not appear in the L^1 penalty, we have

$$\hat{\beta}_0 = (\mathbf{1}^\top \mathbf{1})^{-1} \mathbf{1}^\top Y = \bar{Y}.$$

Thus, the LASSO problem reduces to

$$\min_{\beta_0 \in \mathbb{R}, \gamma \in \mathbb{R}^{p-1}} \{ \|Y - \beta_0 \mathbf{1} - Z\gamma\|_2^2 + \lambda \|\gamma\|_1 \} = \min_{\gamma \in \mathbb{R}^{p-1}} \{ \|u - Z\gamma\|_2^2 + \lambda \|\gamma\|_1 \}.$$

where $u = Y - \bar{Y}\mathbf{1}$ for tidiness. Expanding $\|u - Z\gamma\|_2^2$ gives

$$u^\top u - 2u^\top Z\gamma + \gamma^\top \underbrace{(Z^\top Z)\gamma}_{=I} = u^\top u - 2\underbrace{Y^\top Z\gamma}_{=\hat{\gamma}^\top} + 2\bar{Y}\mathbf{1}^\top \underbrace{Z\gamma}_{=0} + \gamma^\top \gamma$$

Since $u^\top u$ does not depend on γ , we see that the LASSO objective function is

$$-2\hat{\gamma}^\top \gamma + \|\gamma\|_2^2 + \lambda \|\gamma\|_1.$$

Clearly, this has the same minimizer if multiplied across by $1/2$, i.e.

$$-\hat{\gamma}^\top \gamma + \frac{1}{2} \|\gamma\|_2^2 + \frac{1}{2} \lambda \|\gamma\|_1 = \sum_{i=1}^{p-1} (-\hat{\gamma}_i \gamma_i + \frac{1}{2} \gamma_i^2 + \frac{\lambda}{2} |\gamma_i|).$$

Shrinkage XII

Notice that we now have a sum of $p - 1$ objective functions, each depending only on one γ_i . We can thus optimise each separately. That is, for any given $i \leq p - 1$, we must minimise

$$-\hat{\gamma}_i \gamma_i + \frac{1}{2} \gamma_i^2 + \frac{\lambda}{2} |\gamma_i|.$$

We distinguish 3 cases:

- ① **Case $\hat{\gamma}_i = 0$.** In this case, the objective function becomes $\frac{1}{2} \gamma_i^2 + \frac{\lambda}{2} |\gamma_i|$ and it is clear that it is minimised when $\gamma_i = 0$. So in this case $\check{\gamma}_i = 0$.
- ② **Case $\hat{\gamma}_i > 0$.** In this case, the objective function $-\hat{\gamma}_i \gamma_i + \frac{1}{2} \gamma_i^2 + \frac{\lambda}{2} |\gamma_i|$ is minimised somewhere in the range $\gamma_i \in [0, \infty)$ because the term $-\hat{\gamma}_i \gamma_i$ is negative there (and all other terms are positive). But when $\gamma_i \geq 0$, the objective function becomes

$$-\hat{\gamma}_i \gamma_i + \frac{1}{2} \gamma_i^2 + \frac{\lambda}{2} \gamma_i = \left(\frac{\lambda}{2} - \hat{\gamma}_i \right) \gamma_i + \frac{1}{2} \gamma_i^2.$$

If $\frac{\lambda}{2} - \hat{\gamma}_i \geq 0$, then the minimum over $\gamma_i \in [0, \infty)$ is clearly at $\gamma_i = 0$.

Otherwise, when $\frac{\lambda}{2} - \hat{\gamma}_i < 0$, we differentiate and find the minimum at

$\gamma_i = \hat{\gamma}_i - \lambda/2 > 0$. In summary, $\check{\gamma}_i = (\hat{\gamma}_i - \lambda/2)_+ = \min(\hat{\gamma}_i, (\hat{\gamma}_i - \lambda/2)_+)$.

Shrinkage XIII

- ③ Case $\hat{\gamma}_i < 0$. In this case, the objective function $-\hat{\gamma}_i \gamma_i + \frac{1}{2} \gamma_i^2 + \frac{\lambda}{2} |\gamma_i|$ is minimised somewhere in the range $\gamma_i \in (-\infty, 0]$ because the term $-\hat{\gamma}_i \gamma_i$ is negative there (and all other terms are positive). But when $\gamma_i \leq 0$, the objective function becomes

$$-\hat{\gamma}_i \gamma_i + \frac{1}{2} \gamma_i^2 + \frac{\lambda}{2} (-\gamma_i) = \left(\frac{\lambda}{2} + \hat{\gamma}_i \right) (-\gamma_i) + \frac{1}{2} \gamma_i^2 = \left(\frac{\lambda}{2} - |\hat{\gamma}_i| \right) (-\gamma_i) + \frac{1}{2} \gamma_i^2.$$

If $\frac{\lambda}{2} - |\hat{\gamma}_i| \geq 0$, then the minimum over $\gamma_i \in (-\infty, 0]$ is clearly at $\gamma_i = 0$, since $-\gamma_i$ ranges over $[0, \infty)$. Otherwise, when $\frac{\lambda}{2} - |\hat{\gamma}_i| < 0$, we differentiate and find the minimum at $\gamma_i = -|\hat{\gamma}_i| + \lambda/2 < 0$, which we may re-write as:

$$-|\hat{\gamma}_i| + \lambda/2 = -(|\hat{\gamma}_i| - \lambda/2) = \text{sgn}(\hat{\gamma}_i) (|\hat{\gamma}_i| - \lambda/2).$$

In summary, $\check{\gamma}_i = \text{sgn}(\hat{\gamma}_i) (|\hat{\gamma}_i| - \lambda/2)_+$.

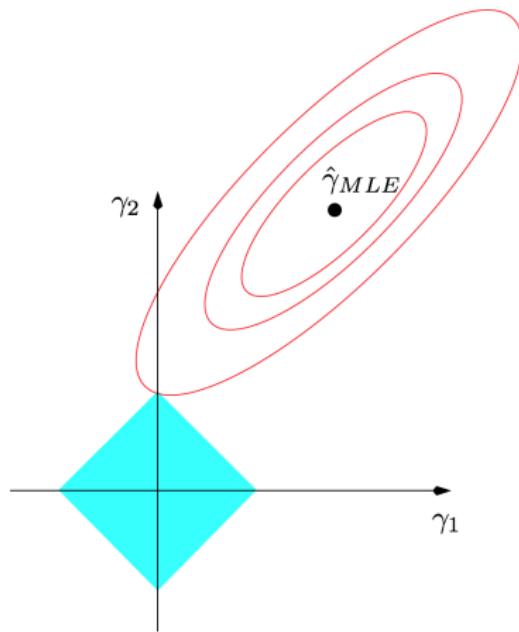
The proof is now complete, as we can see that all three cases yield

$$\check{\gamma}_i = \text{sgn}(\hat{\gamma}_i) \left(|\hat{\gamma}_i| - \frac{\lambda}{2} \right)_+, \quad i = 1, \dots, p-1.$$



Shrinkage XIV

---, ---

Figure: L^1 Shrinkage (the LASSO)

Shrinkage XV

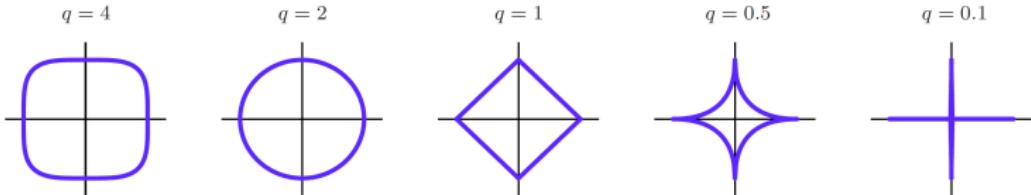
Intuition: L_1 norm induces “sharp” balls!

- Balls more concentrated around the axes
- Induces model selection by regulating the LASSO (through λ)

Extreme case: L^0 “Norm”, gives best subset selection!

$$\|\gamma\|_0 = \sum_{j=1}^{p-1} |\gamma_j|^0 = \sum_{j=1}^{p-1} \mathbf{1}_{\{\gamma_j \neq 0\}} = \#\{j : \gamma_j \neq 0\}$$

Generally: $\|\gamma\|_q^q = \sum_{j=1}^{p-1} |\gamma_j|^q$, sharp balls for $0 < q \leq 1$



But L^1 gives sharpest **convex** ball among these.

Shrinkage XVI

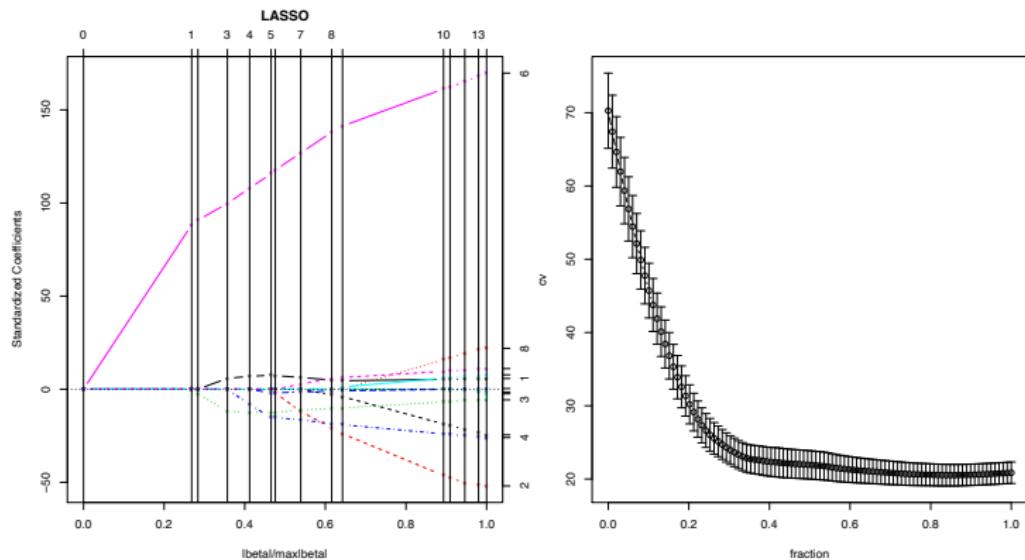


Figure: LASSO and CV for different values of $r(\lambda)/\|\hat{\gamma}\|_1$ for the bodyfat data (LARS algorithm)

Generalised Linear Model

Back to the big picture:

$$Y \text{ (random output)} \xleftarrow{\text{whose law is influenced by}} x \text{ (non-random input)}$$

General formulation:

$$Y_i \stackrel{\text{independent}}{\sim} \text{Distribution} \left\{ g(x_i) \right\}, \quad i = 1, \dots, n.$$

$$= \theta_i$$

Distribution / Function g	$g(x_i^\top) = x_i^\top \beta$	g nonparametric
Gaussian	Linear Regression ✓	Smoothing
Exponential Family	GLM ←	GAM

Generalised Linear Model II

Generalised Linear Models: regression with exponential family responses!

GLM for Y_1, \dots, Y_n

Response vector $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ has **independent entries** with joint law

$$f(\mathbf{y}; \boldsymbol{\phi}) = \prod_{i=1}^n \exp\{\phi_i y_i - \gamma(\phi_i) + S(y_i)\} = \exp\left\{\boldsymbol{\phi}^\top \mathbf{y} - \sum_{i=1}^n \gamma(\phi_i) + \sum_{i=1}^n S(y_i)\right\},$$

where $\boldsymbol{\phi} \in \Phi \subseteq \mathbb{R}^n$ is the **natural parameter** with Φ open. The parameter varies as a function of the covariates via

$$\boldsymbol{\phi} = \mathbf{X}_n \boldsymbol{\beta},$$

for \mathbf{X}_n the $n \times p$ covariate matrix of rank p and $\boldsymbol{\beta}$ a p -dimensional parameter.

In our general notation

$$Y_i \stackrel{\text{independent}}{\sim} f(y_i; \phi_i) = \exp\{\underbrace{(\mathbf{x}_i^\top \boldsymbol{\beta})}_{=\phi_i} y_i - \underbrace{\gamma(\mathbf{x}_i^\top \boldsymbol{\beta})}_{=\phi_i} + S(y_i)\}, \quad i = 1, \dots, n.$$

where the row vector \mathbf{x}_i^\top is the i th row of \mathbf{X}_n .

Generalised Linear Model III

467 / 561

Comments:

- Notice that the sufficient statistic for each marginal distribution $f(y; \phi_i)$ was taken to be the identity $T(Y) = Y$.
- This does not incur any loss of generality for two reasons:
 - ① In the three main GLM of interest (Gaussian, Bernoulli, and Poisson) the natural statistic is for $f(y; \phi_i)$ is indeed the identity.
 - ② More generally, since we only observe a single observation from each $f(y; \phi_i)$, if the natural statistic were $T(Y_i) \neq Y_i$, we could re-define the response to just be $T_i = T(Y_i)$. The sampling distribution of T_i can be shown to also be a one-parameter exponential family with the same natural parameter.
- Recall from our sampling theory results:
 - $\mu_i = \mathbb{E}[Y_i] = \frac{\partial}{\partial \phi_i} \gamma(\phi_i)$
 - $\text{var}[Y_i] = \frac{\partial^2}{\partial \phi_i^2} \gamma(\phi_i)$
 - So if γ is invertible (which it is for the three main examples), the variance can be written as a function of the mean:

$$\text{var}[Y_i] = \gamma''([\gamma']^{-1}(\mu_i)) = V(\mu_i).$$

SGeneralised Linear Model IV

Interpretation of $\phi_i = \mathbf{x}_i^\top \boldsymbol{\beta}$?

- In key cases ϕ_i is directly interpretable. If not, can switch perspective using the mean μ_i as defining parameter, connected to the linear predictor $\mathbf{x}_i^\top \boldsymbol{\beta}$ via

$$[\gamma']^{-1}(\mu_i) = \mathbf{x}_i^\top \boldsymbol{\beta} = \phi_i$$

- The function $[\gamma']^{-1}(\cdot)$ is called the natural link function.
- Instead of $[\gamma']^{-1}$ can use other link functions $g(\cdot)$ and postulate

$$g(\mu_i) = \mathbf{x}_i^\top \boldsymbol{\beta}.$$

This will also yield a GLM, but now the natural parameter will not be equal to the linear predictor but to some function $u(\mathbf{x}_i^\top \boldsymbol{\beta})$ of it.

- In summary, the nomenclature is:
 - $g(\cdot)$ is the link function
 - $h = g^{-1}$ is the inverse link function
 - $g(\cdot) = [\gamma']^{-1}(\cdot)$ is the natural link function
- Will focus on natural links but methods/results generalise quite easily.

Generalised Linear Model V

With natural link, the **loglikelihood** (up to constants w.r.t. β) is

$$\ell_n(\beta) = \beta^\top X_n^\top Y - \sum_{i=1}^n \gamma(x_i^\top \beta)$$

for x_i^\top the i th row of X_n . The corresponding $p \times 1$ **derivative (score function)** is

$$\nabla_\beta \ell_n(\beta) = X_n^\top Y - \sum_{i=1}^n x_i \gamma'(x_i^\top \beta) = \sum_{i=1}^n x_i (Y_i - \mu_i) = X_n^\top (Y - \mu)$$

with $p \times p$ **covariance equaling the information matrix** and given by

$$\begin{aligned} \text{cov}\{\nabla_\beta \ell_n(\beta)\} &= \sum_{i=1}^n x_i^\top \text{cov}\{Y_i - \mu_i\} x_i = X_n^\top V(\beta) X_n \\ &= -\nabla_\beta^2 \ell_n(\beta) = I_n(\beta), \end{aligned}$$

where $\text{cov}\{Y\} = V(\beta) \succ 0$ is diagonal, with i th diagonal element

$$\text{var}\{Y_i\} = \gamma''(\phi_i) = \gamma''(x_i^\top \beta).$$

Generalised Linear Model VI

Thus, if the MLE $\hat{\beta}$ exists it is also unique, and must satisfy:

$$\sum_{i=1}^n \mathbf{x}_i \left(Y_i - \gamma'(\mathbf{x}_i^\top \hat{\beta}) \right) = 0$$

By a first order Taylor expansion of γ' , we have

$$\gamma'(\mathbf{x}_i^\top \hat{\beta}) \approx \gamma'(\mathbf{x}_i^\top \tilde{\beta}) + \mathbf{x}_i^\top (\tilde{\beta} - \hat{\beta}) \gamma''(\mathbf{x}_i^\top \tilde{\beta})$$

for some guesstimate $\tilde{\beta}$ near $\hat{\beta}$. Plugging into the score equation yields

$$\sum_{i=1}^n \mathbf{x}_i^\top \left(Y_i - \gamma'(\mathbf{x}_i^\top \tilde{\beta}) + \mathbf{x}_i^\top (\tilde{\beta} - \hat{\beta}) \gamma''(\mathbf{x}_i^\top \tilde{\beta}) \right) \approx 0$$

$$\implies \sum_{i=1}^n \gamma''(\mathbf{x}_i^\top \tilde{\beta}) \mathbf{x}_i^\top (Z_i - \mathbf{x}_i^\top \hat{\beta}) = \mathbf{X}_n^\top \mathbf{V}(\tilde{\beta})(\tilde{\mathbf{Z}} - \mathbf{X}_n \hat{\beta}) \approx 0$$

where we defined $\tilde{\mathbf{Z}} = (\tilde{Z}_1, \dots, \tilde{Z}_n)^\top$ to be the adjusted response

$$\tilde{Z}_i = \mathbf{x}_i^\top \tilde{\beta} + \frac{1}{\gamma''(\mathbf{x}_i^\top \tilde{\beta})} (Y_i - \gamma'(\mathbf{x}_i^\top \tilde{\beta})) \quad \text{so} \quad \tilde{\mathbf{Z}} = \mathbf{X}_n \tilde{\beta} + \mathbf{V}^{-1}(\tilde{\beta})(\mathbf{Y} - \mu(\tilde{\beta})).$$

Generalised Linear Model VII

The last expression for the score expression now yields:

$$\hat{\beta} \approx (X_n^\top V(\tilde{\beta}) X_n)^{-1} X_n^\top V(\tilde{\beta}) \tilde{Z}$$

Just a **weighted least squares estimate!** Where's the catch?

- Weight matrix $V(\tilde{\beta})$ requires specification of an initial guesstimate $x_i^\top \tilde{\beta}$ sufficiently close to $x_i^\top \hat{\beta}$.
- Luckily we can give such a guesstimate by recalling that $\mu_i = \gamma'(x_i^\top \beta)$ so that estimating μ_i by Y_i yields the guesstimate $x_i^\top \tilde{\beta} \equiv (\gamma')^{-1}(Y_i)$.
- Suggests the following **Iteratively Reweighted Least Squares (IRLS)**

IRLS

① Initialize with $x_i^\top \beta^{(0)} \leftarrow \gamma'(Y_i)$ and $Z_i^{(0)} = x_i^\top \beta^{(0)} + \frac{(Y_i - \gamma'(x_i^\top \beta^{(0)}))}{\gamma''(x_i^\top \beta^{(0)})}$

② Update with $\beta^{(j+1)} \leftarrow (X_n^\top V(\beta^{(j)}) X_n)^{-1} X_n^\top V(\beta^{(j)}) Z^{(j)}$

- Equivalent to Newton-Raphson iteration.
- Not always guaranteed to converge.

Generalised Linear Model VIII

(Approximate) Sampling Distribution of MLE in GLM

Heuristics:

- Suppose we had started iteration at true β and stopped at first iterate:

$$\hat{\beta} = (X_n^\top V(\beta) X_n)^{-1} X_n^\top V(\beta) Z$$

where

$$Z_i = x_i^\top \beta + \frac{1}{\gamma''(x_i^\top \beta)} (Y_i - \gamma'(x_i^\top \beta)) \quad \text{so} \quad Z = X_n \beta + V^{-1}(\beta)(Y - \mu(\beta))$$

- This would give us,

$$\hat{\beta} = \beta + (X_n^\top V(\beta) X_n)^{-1} X_n^\top (Y - \mu(\beta))$$

- So we would expect $\mathbb{E}[\hat{\beta}] = \beta$ and $\text{cov}[\hat{\beta}] = (X_n^\top V^{-1}(\beta) X_n)^{-1} = I_n^{-1}(\beta)$.
- And we would conjecture a Gaussian limiting law (paralleling the IID setting)

Under conditions, this is indeed what we obtain.

Result stated in form valid for more general (sufficiently regular) link functions.

Generalised Linear Model IX

Theorem (Asymptotic Normality of MLE in GLM)

In the same context and notation as before, assume that:

(C1) $\beta \in B$ for B an open convex subset of \mathbb{R}^p .

(C2) The $p \times p$ matrix $X_n^\top X_n$ is of full rank for all n .

(C3) The information diverges, i.e. $\lambda_{\min}(\mathcal{I}_n(\beta)) \rightarrow \infty$ as $n \rightarrow \infty$ for $\lambda_{\min}(\cdot)$ the smallest eigenvalue.

(C4) Given any parameter $\beta \in \mathbb{R}^p$ it holds that

$$\sup_{\alpha \in N_\delta(\beta)} \left\| \mathcal{I}_n^{-1/2}(\beta) \mathcal{I}_n^{1/2}(\alpha) - I_{p \times p} \right\| \rightarrow 0$$

$\forall \delta > 0$, where $N_\delta(\beta) = \{\alpha \in \mathbb{R}^p : (\alpha - \beta)^\top \mathcal{I}_n(\beta)(\alpha - \beta) \leq \delta\}$.

Then, as $n \rightarrow \infty$, provided it exists, the MLE $\hat{\beta}_n$ of β_0 is unique & satisfies

$$\mathcal{I}_n^{1/2}(\beta_0)(\hat{\beta}_n - \beta_0) \xrightarrow{d} N(0, I_{p \times p}).$$

- Recall that under canonical link $\mathcal{I}_n(\beta) = X_n^\top V(\beta) X_n$.