

# Statistics for Data Science: Week 12

Myrto Limnios and Rajita Chandak

Institute of Mathematics – EPFL

`rajita.chandak@epfl.ch, myrto.limnios@epfl.ch`



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

Multicollinearity **problem** is that  $\det [\mathbf{X}^\top \mathbf{X}] \approx 0$   
[i.e.  $\mathbf{X}^\top \mathbf{X}$  *almost not invertible*]

**A Solution:** add a “small amount” of a full rank matrix to  $\mathbf{X}^\top \mathbf{X}$ .

For reasons to become clear soon, we *standardise* the design matrix:

- Write  $\mathbf{X} = (\mathbf{1} \ \mathbf{W})$ ,  $\beta = (\beta_0 \ \gamma)^\top$
- Recentre/rescale the covariates (columns) defining:  $\mathbf{Z}_j = \frac{\sqrt{n}}{\text{sd}(\mathbf{W}_j)} (\mathbf{W}_j - \mathbf{1} \overline{\mathbf{W}_j})$ 
  - Coefficients now have common scale
  - Interpretation of  $\beta_j$  slightly different: not “mean impact on response per unit change of explanatory variable”, but now “mean impact on response per unit deviation of explanatory variable from its mean, measured in units of standard deviation”
- The  $\mathbf{Z}_j$  are all orthogonal to  $\mathbf{1}$  and are of unit norm.

- Since  $\mathbf{Z}_j \perp \mathbf{1}$  for all,  $j$ , we can estimate  $\beta_0$  and  $\gamma$  by two separate regressions (orthogonality).
- Least squares estimators based on the *standardized design matrix* become

$$\hat{\beta}_0 = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i, \quad \hat{\gamma} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y}.$$

- Ridge regression replaces  $\mathbf{Z}^\top \mathbf{Z}$  by  $\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_{(p-1) \times (p-1)}$  (i.e. adds a “ridge”)

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

Adding  $\lambda \mathbf{I}_{(p-1) \times (p-1)}$  to  $\mathbf{Z}^\top \mathbf{Z}$  makes inversion more stable  
 $\hookrightarrow \lambda$  called *ridge parameter*.

→ Ridge term  $\lambda I$  seems slightly ad-hoc. Motivation?

↪ Can see that  $(\hat{\beta}_0 \quad \hat{\gamma}) = (\bar{Y} \quad (\mathbf{Z}^\top \mathbf{Z} + \lambda I)^{-1} \mathbf{Z}^\top \mathbf{Y})$  minimizes

$$\|\mathbf{Y} - \beta_0 \mathbf{1} - \mathbf{Z}\gamma\|_2^2 + \lambda \|\gamma\|_2^2$$

or equivalently

$$\|\mathbf{Y} - \beta_0 \mathbf{1} - \mathbf{Z}\gamma\|_2^2 \quad \text{subject to} \quad \sum_{j=1}^{p-1} \gamma_j^2 = \|\gamma\|_2^2 \leq r(\lambda)$$

instead of least squares estimator which minimizes

$$\|\mathbf{Y} - \beta_0 \mathbf{1} - \mathbf{Z}\gamma\|_2^2.$$

Idea: in the presence of collinearity, coefficients are ill-defined: a wildly positive coefficient can be cancelled out by a largely negative coefficient (many coefficient combinations can produce the same effect). By imposing a *size* constraint, we limit the possible coefficient combinations!

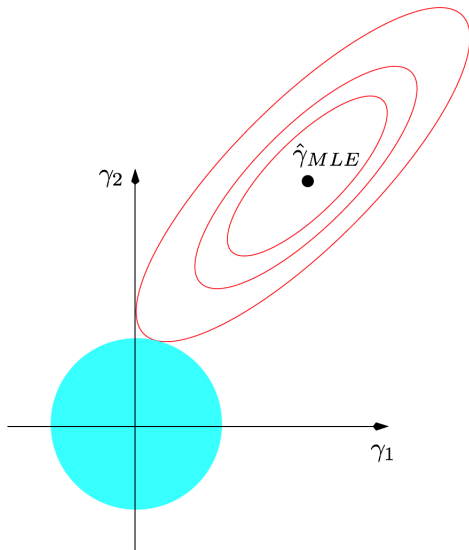


Figure:  $L^2$  Shrinkage [Ridge Regression]

## Theorem

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

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

is

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

## Proof.

Write

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

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

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

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

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

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

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

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

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

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

$$\hat{\gamma} = \sum_{j=1}^q \frac{\lambda_j}{\lambda_j^2 + \lambda} (\mathbf{v}_j^\top \mathbf{Y}) \mathbf{u}_j,$$

where the  $\mathbf{V}_j$ s and  $\mathbf{U}_j$ s are the columns of  $\mathbf{V}$  and  $\mathbf{U}$ , respectively.

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

$$\hat{\gamma} = \sum_{j=1}^q \frac{1}{\lambda_j} (\mathbf{v}_j^\top \mathbf{Y}) \mathbf{u}_j,$$

which is not even defined if  $\mathbf{Z}$  is of reduced rank.

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



## Proposition

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

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

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. □

## Corollary (Domination over Least Squares)

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

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

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  $\lambda$  depends on unknowns.
- Choosing a good  $\lambda$  is all about balancing bias and variance.

## Proof.

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

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

To go from 2nd to 3rd line, we wrote

$$\begin{aligned} \sigma^2(\mathbf{Z}^\top \mathbf{Z})^{-1} &= \sigma^2(\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}) (\mathbf{Z}^\top \mathbf{Z})^{-1} (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}) (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} \\ &= (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} (\sigma^2 \mathbf{Z}^\top \mathbf{Z} + 2\sigma^2 \lambda \mathbf{I} + \sigma^2 \lambda^2 (\mathbf{Z}^\top \mathbf{Z})^{-1}) (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I})^{-1} \end{aligned}$$

and did the tedious (but straightforward) algebra. The **final 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

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

for  $\{\gamma/\|\gamma\|, \theta_1, \dots, \theta_{q-1}\}$  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.



## Bias–Variance Tradeoff, Again.

Role of  $\lambda$ : Regulates Bias–Variance tradeoff

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

Recall:

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

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

So choose  $\lambda$  so as to optimally to balance **bias**/**variance**

Use cross validation!



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!

When the explanatory variables are orthogonal (i.e.  $\mathbf{Z}^\top \mathbf{Z} = \mathbf{I}$ ), then the LASSO<sup>1</sup> exactly performs model selection via thresholding:

## Theorem (Orthogonal Design $\leftrightarrow$ Model Selection)

*Consider the linear model*

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

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

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

*Then, the unique solution to the LASSO problem*

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

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

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

---

<sup>1</sup>Least Absolute Shrinkage and Selection Operator.

## Proof (\*).

Note that since  $Z^\top \mathbf{1} = 0$  and since  $\beta_0$  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

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

Since  $u^\top u$  does not depend on  $\gamma$ , 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.

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

Notice that we now have a sum of  $p - 1$  objective functions, each depending only on one  $\gamma_i$ . 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:

- ❶ **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  $\check{\gamma}_i = 0$ .**
- ❷ **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,  $\check{\gamma}_i = (\hat{\gamma}_i - \lambda/2)_+ = \text{sgn}(\hat{\gamma}_i)(|\hat{\gamma}_i| - \lambda/2)_+$ .**



- 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

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



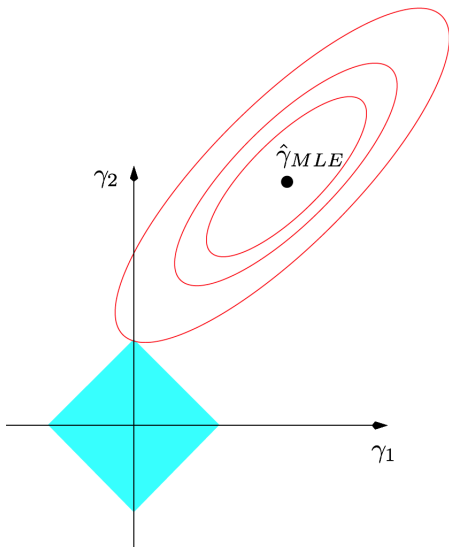


Figure:  $L^1$  Shrinkage (the LASSO)

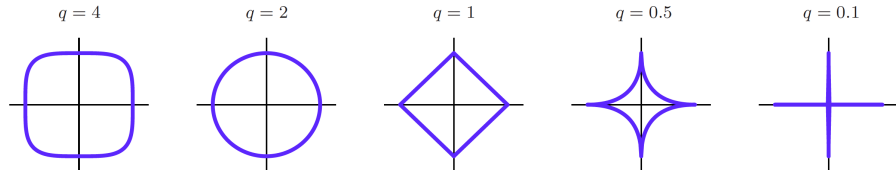
**Intuition:**  $L_1$  norm induces “sharp” balls!

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

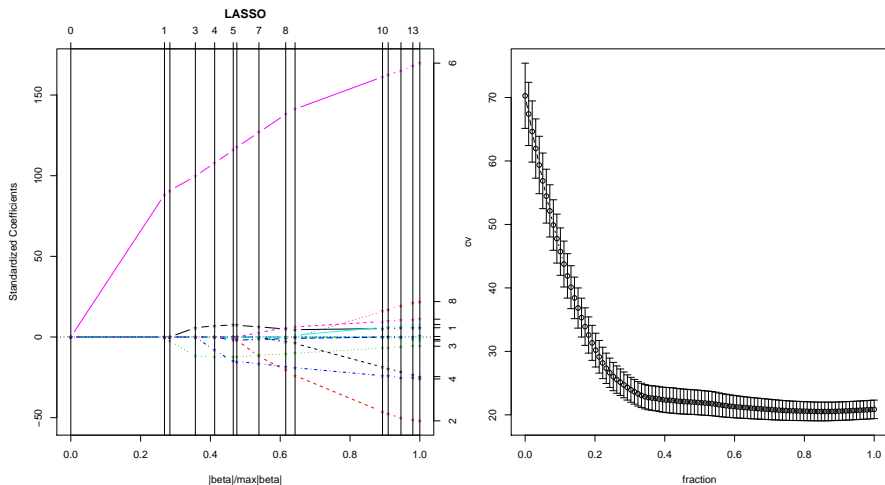
Extreme case:  $L^0$  “Norm”, gives best subset selection!

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

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.



**Figure:** LASSO and CV for different values of  $r(\lambda)/\|\hat{y}\|_1$  for the bodyfat data (LARS algorithm), on the right: fraction of the final  $L^1$  norm

- Regularization achieves smaller variance at the price of non-zero bias
- Need to find the best trade-off between variance/bias
- The larger the parameter  $\lambda$ , the greater the shrink (so different values of  $\lambda$  imply different estimates for  $\gamma$ )
- Ridge regression is a convex objective risk, and works well when there is a subset of coefficients that are small
- LASSO is used as *feature selection* method, as it yields estimates of parameters to be exactly equal to 0
- The tuning parameter  $\lambda$  tailors the strength of the penalization
- Choosing numerically the tuning parameter is often by  $k$ -fold cross validation (and NOT based on training error): we want to avoid over-/under-fitting

# Generalised Linear Models

Back to the big picture:

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

General formulation:

$$Y_i \overset{\text{independent}}{\sim} \underbrace{\text{Distribution}\{g(x_i)\}}_{=\theta_i}, \quad i = 1, \dots, n.$$

Distribution / Function $g$	$g(\mathbf{x}_i^\top) = \mathbf{x}_i^\top \boldsymbol{\beta}$	$g$ nonparametric
Gaussian	Linear Regression ✓	Smoothing
Exponential Family	GLM $\leftarrow$	GAM

## 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

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

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

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

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

In our general notation

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

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



## 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  $\gamma$  is invertible (which it is for the three main examples), the variance can be written as a function of the mean:

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

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

- In key cases  $\phi_i$  is directly interpretable. If not, can switch perspective using the mean  $\mu_i$  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.

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

$$\ell_n(\beta) = \beta^\top \mathbf{X}_n^\top \mathbf{Y} - \sum_{i=1}^n \gamma(\mathbf{x}_i^\top \beta)$$

for  $\mathbf{x}_i^\top$  the  $i$ th row of  $\mathbf{X}_n$ . The corresponding  $p \times 1$  **derivative (score function)** is

$$\nabla_{\beta} \ell_n(\beta) = \mathbf{X}_n^\top \mathbf{Y} - \sum_{i=1}^n \mathbf{x}_i \gamma'(\mathbf{x}_i^\top \beta) = \sum_{i=1}^n \mathbf{x}_i (Y_i - \mu_i) = \mathbf{X}_n^\top (\mathbf{Y} - \boldsymbol{\mu})$$

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

$$\begin{aligned} \text{cov} \{ \nabla_{\beta} \ell_n(\beta) \} &= \sum_{i=1}^n \mathbf{x}_i^\top \text{cov} \{ Y_i - \mu_i \} \mathbf{x}_i = \mathbf{X}_n^\top \mathbf{V}(\beta) \mathbf{X}_n \\ &= -\nabla_{\beta}^2 \ell_n(\beta) = \mathcal{I}_n(\beta), \end{aligned}$$

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

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

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  $\gamma'$ , 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

$$\begin{aligned} \sum_{i=1}^n \mathbf{x}_i^\top \left( Y_i - \gamma'(\mathbf{x}_i^\top \tilde{\beta}) + \mathbf{x}_i^\top (\tilde{\beta} - \hat{\beta}) \gamma''(\mathbf{x}_i^\top \tilde{\beta}) \right) &\approx 0 \\ \implies \sum_{i=1}^n \gamma''(\mathbf{x}_i^\top \tilde{\beta}) \mathbf{x}_i^\top (Z_i - \mathbf{x}_i^\top \hat{\beta}) &= \mathbf{X}_n^\top \mathbf{V}(\tilde{\beta}) (\tilde{\mathbf{Z}} - \mathbf{X}_n \hat{\beta}) \approx 0 \end{aligned}$$

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

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

The last expression for the score expression now yields:

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

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

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

## IRLS

- 1 Initialize with  $\mathbf{x}_i^\top \beta^{(0)} \leftarrow (\gamma')^{-1}(Y_i)$  and  $Z_i^{(0)} = \mathbf{x}_i^\top \beta^{(0)} + \frac{(Y_i - \gamma'(\mathbf{x}_i^\top \beta^{(0)}))}{\gamma''(\mathbf{x}_i^\top \beta^{(0)})}$
- 2 Update with  $\beta^{(j+1)} \leftarrow (\mathbf{X}_n^\top \mathbf{V}(\beta^{(j)}) \mathbf{X}_n)^{-1} \mathbf{X}_n^\top \mathbf{V}(\beta^{(j)}) \mathbf{Z}^{(j)}$ 
  - Equivalent to Newton-Raphson iteration.
  - **Not always guaranteed to converge.**

## Heuristics:

- Suppose we had started iteration at true  $\beta$  and stopped at first iterate:

$$\hat{\beta} = (\mathbf{X}_n^\top \mathbf{V}(\beta) \mathbf{X}_n)^{-1} \mathbf{X}_n^\top \mathbf{V}(\beta) \mathbf{Z}$$

where

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

- This would give us,

$$\hat{\beta} = \beta + (\mathbf{X}_n^\top \mathbf{V}(\beta) \mathbf{X}_n)^{-1} \mathbf{X}_n^\top (\mathbf{Y} - \boldsymbol{\mu}(\beta))$$

- So we would expect  $\mathbb{E}[\hat{\beta}] = \beta$  and  $\text{cov}[\hat{\beta}] = (\mathbf{X}_n^\top \mathbf{V}^{-1}(\beta) \mathbf{X}_n)^{-1} = \mathcal{I}_n^{-1}(\beta)$ .
- 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.

## 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  $\mathbf{X}_n^\top \mathbf{X}_n$  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

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

$\forall \delta > 0$ , 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  $\beta_0$  is unique & satisfies

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

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

## Comments on conditions:

- (C1) implies that  $\mathbf{X}_n\beta$  ranges over an open set, and so  $\gamma$  is infinitely differentiable, and our exponential family possesses all moments.
- (C2) is as with our linear model, and essentially means that our covariates should not be perfectly correlated.
- (C3) is similar to the “balanced design” assumption we had for the asymptotics of a non-Gaussian linear model.
- (C1-C3) are up to us: they depend on the design matrix  $\mathbf{X}_n$ , which in principle is for us to choose.
- (C4) asks that the (root) information matrix converge uniformly on compact ellipsoids centred at the true parameter.



### Comments on conclusion:

- Can also be read as saying that for  $n$  sufficiently large,

$$\hat{\beta}_n \stackrel{d}{\approx} N(\beta_0, \mathcal{I}_n^{-1}(\beta_0))$$

- Conclusion also immediately implies that

$$(\hat{\beta}_n - \beta_0)^\top \mathcal{I}_n(\beta_0) (\hat{\beta}_n - \beta_0) \stackrel{d}{\rightarrow} \chi_p^2.$$

- Allows adapting testing/CI developed before using LR and Wald statistics.

In Gaussian linear regression we used **sums of squares** to measure fit and compare nested models. **What about in GLM?**

**Idea:** compare best possible to observed maximised loglikelihood

- Let  $\ell_n(\hat{\beta}) = \hat{\beta}^\top \mathbf{X}^\top \mathbf{Y} + \sum_{i=1}^n \gamma(\mathbf{x}_i^\top \hat{\beta})$  be the maximised loglikelihood.
- Define the **saturated model** to be that which has

$$\# \text{parameters} = \# \text{observations}$$

i.e. where we replace  $\mathbf{X}_n \beta$  by some unconstrained  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n)^\top \in \mathbb{R}^n$ .  
(and thus we also replace  $\mathbf{x}_i^\top \beta$  by  $\eta_i$ ).

- Let  $\hat{\boldsymbol{\eta}}$  be the maximiser of  $\ell_n(\boldsymbol{\eta}) = \boldsymbol{\eta}^\top \mathbf{Y} + \sum_{i=1}^n \gamma(\eta_i)$  w.r.t.  $\boldsymbol{\eta}$ .
- Define the **saturated loglikelihood** as  $\ell_n(\hat{\boldsymbol{\eta}})$ .

## Definition ((Scaled) Deviance)

$$D = 2(\ell_n(\hat{\boldsymbol{\eta}}) - \ell_n(\hat{\boldsymbol{\beta}})) = 2\left((\hat{\boldsymbol{\eta}} - \mathbf{X}_n\hat{\boldsymbol{\beta}})^\top \boldsymbol{\Upsilon} + \sum_{i=1}^n (\gamma(\hat{\eta}_i) - \gamma(\mathbf{x}_i^\top \hat{\boldsymbol{\beta}}))\right)$$

### Comments:

- Always  $D \geq 0$  (why?)
- Small  $D$  implies a good model fit ( $\mathbf{X}_n\hat{\boldsymbol{\beta}} \approx \hat{\boldsymbol{\eta}}$ ).
- Large  $D$  implies poor fit.
- In Gaussian case: deviance  $\equiv$  residual sum of squares (exercise).
- Can now use deviance differences to mimic sum-of-square ratios and construct a GLM variant of the F-test.

- Consider the problem of comparing two nested models:
  - Model A:**  $\beta = (\beta_1, \dots, \beta_p)^\top \in \mathbb{R}^p$  vary freely — MLE  $\hat{\beta}^A$
  - Model B:** for  $q < p$ ,  $(\beta_1, \dots, \beta_q) \in \mathbb{R}^q$  vary freely, but  $\beta_{q+1}, \dots, \beta_p$  are fixed — hence only  $q$  free parameters, with MLE  $\hat{\beta}^B$
- Model  $B$  is *nested within* model  $A$ :  $B$  can be obtained by restrictions on  $A$ 
  - More generally, could have Model  $B$  with  $\beta$  constrained to vary in a subspace  $\mathcal{V}$  of dimension  $q < p$ , which we can write as  $\beta = \underbrace{\mathbf{Q}_{p \times q}}_{q \times 1} \zeta$  for  $\mathcal{M}(\mathbf{Q}) = \mathcal{V}$ .
- Likelihood ratio test statistic for comparing the models is

$$2(\ell_n(\hat{\beta}^A) - \ell_n(\hat{\beta}^B)) = D_B - D_A,$$

and **when model  $B$  is correct**  $D_B - D_A \xrightarrow{d} \chi_{p-q}^2$ .

**Main idea:** use deviance instead of sums-of-squares and use final iterate of IWLS to get hat matrix

- **Leverage**  $h_{jj}$  defined as  $j$ th diagonal element of

$$\mathbf{H} = \mathbf{V}^{1/2}(\hat{\beta})\mathbf{X}_n(\mathbf{X}_n^\top \mathbf{V}(\hat{\beta})\mathbf{X}_n)^{-1}\mathbf{X}_n^\top \mathbf{V}^{1/2}(\hat{\beta}),$$

- **Cook statistic** now becomes the change in deviance

$$2p^{-1} \left\{ \ell(\hat{\beta}) - \ell(\hat{\beta}_{-j}) \right\},$$

where  $\hat{\beta}_{-j}$  is MLE when  $j$ th case  $(\mathbf{x}_j^\top, Y_j)$  is dropped.

- Cook statistic can be approximated by

$$C_j = \frac{h_{jj}}{p(1 - h_{jj})} r_{Pj}^2,$$

where  $r_{Pj}$  is standardised Pearson residual (to be defined in next slide).

- Deviance residual:

$$d_j = \text{sgn}(\hat{\eta}_j - \mathbf{x}_j^\top \hat{\boldsymbol{\beta}}_j) \left[ 2 \left\{ \underbrace{\eta_j Y_j + \gamma(\eta_j)}_{\ell_j(\hat{\boldsymbol{\eta}})} - \underbrace{[(\mathbf{x}_j^\top \boldsymbol{\beta}) Y_j + \gamma(\mathbf{x}_j^\top \boldsymbol{\beta})]}_{\ell_j(\hat{\boldsymbol{\beta}})} \right\} \right]^{1/2},$$

for which we note that

$$\sum_{j=1}^n d_j^2 = D$$

gives the deviance (in analogy with RSS in Gaussian linear regression).

- Pearson residual:

$$p_j = \frac{Y_i - \gamma'(\mathbf{x}_i^\top \hat{\boldsymbol{\beta}})}{\sqrt{\gamma''(\mathbf{x}_i^\top \hat{\boldsymbol{\beta}})}} = \frac{Y_i - \hat{\mu}_i}{\sqrt{V(\hat{\mu}_i)}} \quad \text{so } \mathbf{r}_P = \mathbf{V}^{-1/2}(\hat{\boldsymbol{\beta}})(\mathbf{Y} - \boldsymbol{\mu}(\hat{\boldsymbol{\beta}})).$$

- Standardised versions:

$$r_{Dj} = \frac{d_j}{(1 - h_{jj})^{1/2}} \quad \& \quad r_{Pj} = \frac{p_j}{(1 - h_{jj})^{1/2}}.$$