

Regression Methods

Myrto Linnios

Autumn 2025 - Week 9

Regularisation - Basic Notions

Tall and wide regressions

$$(X^T X)^{-1} \text{ if } p \gg n$$

- ▶ So far we have supposed that we have a **tall regression**:
 - ▶ the number of units n exceeds the number of variables p ,
 - ▶ the design matrix X has rank p .
- ▶ In many 'modern' settings we instead have a **wide regression**:
 - ▶ n and p are comparable, $p > n$, maybe even $p \gg n$;
 - ▶ in genomics, for example (typically) $n = O(10^2, 10^3)$, $p = O(10^5, 10^6)$;
 - ▶ hence $\text{rank}(X) = \min(n, p) = n$.
- ▶ Even tall X may be 'almost singular', making β 'almost inestimable'.
- ▶ Solutions:
 - ▶ subset selection (drop certain columns of X);
 - ▶ seek different good explanations of response variation, not single model;
 - ▶ regularisation (often with prediction in mind).
- ▶ Certain regularisation methods (e.g., lasso) also perform subset selection.

$p \ll n$

Collinearity

- ▶ Columns of X **collinear** if there exists a non-zero $v_{p \times 1}$ such that $Xv = 0$, i.e., $\text{rank}(X) < p$, so there is no unique $\hat{\beta}$ minimising $\|y - X\beta\|^2$.
- ▶ Software deals with this by dropping columns of X , but it may be better to write $X\beta = XC\gamma$, where XC is full rank and γ has a clear interpretation.
- ▶ If X is nearly collinear, its SVD $U_{n \times n} D_{n \times p} V_{p \times p}^T$, with $d_1 \geq \dots \geq d_p \geq 0$, gives

$$\hat{\beta} = (X^T X)^{-1} X^T y = V D_-^T U^T y = \sum_{r=1}^p \underbrace{(u_r^T y / d_r)}_{\text{diag}} v_r$$

so $\hat{\beta}$ is a linear combination of the vectors v_r with coefficients $u_r^T y / d_r$. As $\text{var}(U^T y) = \sigma^2 I_n$,

$$\text{var}(\hat{\beta}) = \sigma^2 V D_-^T D_- V^T = \sigma^2 \sum_{r=1}^p d_r^{-2} v_r v_r^T,$$

i.e., $\hat{\beta}$ is unstable in the directions corresponding to the v_r with small singular values d_r .

- ▶ In numerical analysis, collinearity often measured using **condition number** $(d_1/d_p)^{1/2}$, but its statistical meaning is unclear.

Regularisation

- ▶ Stop $\hat{\beta}$ from fluctuating too wildly in directions with small eigenvalues d_r , by adding a non-negative penalty $p_\lambda(\beta)$ and choosing β to minimise the **penalised sum of squares**

$$\hat{\beta} \in \arg \min_{\beta} \left[\|y - X\beta\|^2 + p_\lambda(\beta) \right] \quad (1)$$

- ▶ The strength of the penalty depends on a positive parameter λ that constrains β more as λ increases.
- ▶ Often $p_\lambda(\beta) = \lambda p(\beta)$, where, for example,
 - ▶ $p(\beta) = \|\beta\|_2^2 = \sum_{r=1}^p \beta_r^2$ gives **ridge regression** (aka Tikhonov regularisation);
 - ▶ $p(\beta) = \|\beta\|_1 = \sum_{r=1}^p |\beta_r|$ gives the **lasso** (aka L_1 regularisation);
 - ▶ $p(\beta) = (1 - \alpha)\|\beta\|_2^2 + \alpha\|\beta\|_1$ for $0 \leq \alpha \leq 1$ gives the **elastic net**;
 - ▶ $p(\beta) = \sum_{g=1}^G p_g^{1/2} \|\beta_g\|_2$, with β_g being $p_g \times 1$ sub-vectors of β , gives the **grouped lasso**, which penalises factors with parameters β_g .
- ▶ It is useful to see regularisation through the lens of Bayesian inference, with the regularising term equivalent to the prior density.

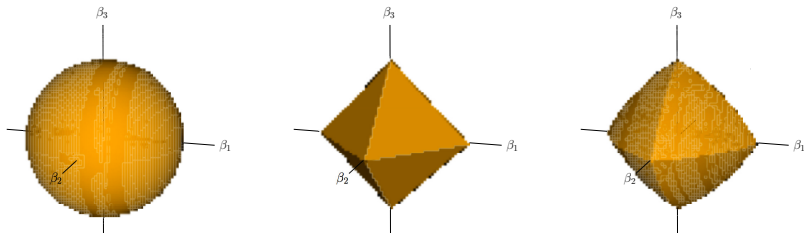
Bound form

- ▶ Equivalently we can take the **bound form** of the minimisation problem, i.e.,

$$\text{minimise}_{\beta} \quad \|y - X\beta\|_2^2 \quad \text{subject to} \quad p(\beta) \leq t,$$

for some $t \geq 0$, where setting $t = \infty$ just gives the least squares estimates.

- ▶ Below: constraint balls for ridge (left), lasso (centre) and elastic-net (right) regularisation. The sharp corners of the last two allow for variable selection as well as shrinkage.



Bayesian approach

- ▶ Treat all unknowns as random variables, and compute conditional distribution of unobserved unknowns conditional on observed unknowns.
- ▶ Requires prior density on β , and if σ^2 is known, then a simple combination of **data model** and **prior model** is

$$y \mid \beta, \sigma^2 \sim \mathcal{N}(X\beta, \sigma^2 I_n), \quad \beta \mid \sigma^2 \sim \mathcal{N}_p(\beta_*, \sigma^2 V_*), \quad (2)$$

where the prior model is determined by β_* and V_* .

- ▶ Full specification would require prior on σ^2 , but we don't need this.
- ▶ Let \equiv mean we have dropped additive constants not involving the argument of a density.
- ▶ The log multivariate normal density is

$$\begin{aligned} \log f(x \mid \mu, \Omega) &= -\frac{m}{2} \log 2\pi - \frac{1}{2} \log |\Omega| - \frac{1}{2} (x - \mu)^T \Omega^{-1} (x - \mu) \\ &\equiv x^T \Omega^{-1} \mu - \frac{1}{2} x^T \Omega^{-1} x \\ &\equiv Q(x) = x^T a - \frac{1}{2} x^T B x, \end{aligned}$$

say, and as $\exp Q(x)$ is proportional to a unique probability density function,

$E(X) = \mu = B^{-1}a$, $\text{var}(X) = \Omega = B^{-1}$, where B is the **precision matrix**.

Bayesian linear model

- ▶ The model (2) gives

$$\begin{aligned}\log f(\beta | y, \sigma^2) &= \log \left\{ \frac{f(y | \beta, \sigma^2) f(\beta | \sigma^2)}{f(y | \sigma^2)} \right\} \\ &\equiv \log f(y | \beta, \sigma^2) + \log f(\beta | \sigma^2) \\ &\equiv -\frac{(y - X\beta)^T (y - X\beta)}{2\sigma^2} - \frac{(\beta - \beta_*)^T V_*^{-1} (\beta - \beta_*)}{2\sigma^2} \\ &\propto \|y - X\beta\|_2^2 + (\beta - \beta_*)^T V_*^{-1} (\beta - \beta_*).\end{aligned}$$

- ▶ Comparison with (1) shows that $p_\lambda(\beta)$ represents prior beliefs about the likely values of β : before seeing the data, the most plausible value is β_* , with precision V_*^{-1} .
- ▶ Dropping more constants,

$$\begin{aligned}\log f(\beta | y, \sigma^2) &\equiv \frac{1}{\sigma^2} \{ \beta^T X^T y - \beta^T (X^T X) \beta / 2 + \beta^T V_*^{-1} \beta_* - \beta^T V_*^{-1} \beta / 2 \} \\ &= \frac{1}{2\sigma^2} \{ 2\beta^T (X^T y + V_*^{-1} \beta_*) - \beta^T (X^T X + V_*^{-1}) \beta \}, \quad (3)\end{aligned}$$

which is $Q(x)$ with x , a and B replaced by β , $(X^T y + V_*^{-1} \beta_*) / \sigma^2$ and $(X^T X + V_*^{-1}) / \sigma^2$.

- ▶ Hence $f(\beta | y, \sigma^2)$ is multivariate normal with mean vector and variance matrix

$$E(\beta | y, \sigma^2) = (X^T X + V_*^{-1})^{-1} (X^T y + V_*^{-1} \beta_*), \quad \text{var}(\beta | y, \sigma^2) = \sigma^2 (X^T X + V_*^{-1})^{-1}.$$

- ▶ The **maximum a posteriori (MAP) estimator** of β is $E(\beta | y, \sigma^2)$, and the MAP estimator of $A_{q \times p} \beta$ is $AE(\beta | y, \sigma^2)$, which has a posterior normal density.
- ▶ When $X^T X$ is invertible,

$$\tilde{\beta} = E(\beta | y, \sigma^2) = (X^T X + V_*^{-1})^{-1}(X^T X \hat{\beta} + V_*^{-1} \beta_*)$$

is an average of $\hat{\beta}$ and β_* , weighted by $X^T X$ and V_*^{-1} .

- ▶ The posterior precision matrix

$$\text{var}(\beta | y, \sigma^2)^{-1} = X^T X / \sigma^2 + V_*^{-1} / \sigma^2$$

adds the Fisher information and the prior precision matrix, V_*^{-1} / σ^2 .

- ▶ High precision corresponds to small variance, and conversely:
 - ▶ letting $V_*^{-1} \rightarrow 0$ yields an improper prior density; and
 - ▶ for large V_*^{-1} the posterior precision is essentially determined by the prior precision.

Thus the prior density regularises $\hat{\beta}$ by including β_* and V_* .

Improper prior density

- ▶ We only need V_* to add information in directions corresponding to small singular values of X , so we might use an **improper prior** in which V_* is singular:

$$f(\beta \mid \sigma^2) = \frac{1}{(2\pi)^{p/2} |V_*|_+^{1/2}} \exp \left\{ -(\beta - \beta_*)^T V_*^{-1} (\beta - \beta_*) / (2\sigma^2) \right\}, \quad (4)$$

where V_* has spectral decomposition ED_*E^T ,

- ▶ $|V_*|_+$ denotes the product of the non-zero elements of D_* , and
 - ▶ $V_*^{-1} = \sum_{r: d_{*r} > 0} e_r e_r^T / d_{*r}$ is a generalized inverse of V_* .
- ▶ Below we write V_*^{-1} even when V_* is invertible.
 - ▶ (4) is improper because it is not integrable in the directions of the columns of E for which the corresponding d_r^* equal zero, but we need only that the posterior density of β be proper, i.e., that the posterior precision matrix

$$\text{var}(\beta \mid y, \sigma^2)^{-1} = X^T X / \sigma^2 + V_*^{-1} / \sigma^2$$

is invertible.

Empirical Bayes

- ▶ Use the data to estimate the prior: construct estimators using Bayesian arguments, but assess their properties using classical criteria (bias, MSE, ...)
- ▶ The estimator $\tilde{\beta} = E(\beta \mid y, \sigma^2)$ has mean and variance

$$\begin{aligned}E(\tilde{\beta} \mid \beta) &= (X^T X + V_*^-)^{-1}(X^T X \beta + V_*^- \beta_*) \\ &= \beta + (X^T X + V_*^-)^{-1} V_*^- (\beta_* - \beta), \\ \text{var}(\tilde{\beta} \mid \beta) &= \sigma^2 (X^T X + V_*^-)^{-1} X^T X (X^T X + V_*^-)^{-1}.\end{aligned}\quad (5)$$

- ▶ Hence $\tilde{\beta}$
 - ▶ is biased unless $\beta_* = \beta$,
 - ▶ has smaller variance than $\hat{\beta}$,

so maybe there is a bias-variance tradeoff when estimating $A\beta$.

- ▶ If we write $\mu = E(\tilde{\beta} | \beta)$, then the MSE is

$$\begin{aligned} E\left(\|A\tilde{\beta} - A\beta\|^2 | \beta\right) &= E\{(\tilde{\beta} - \beta)^T A^T A(\tilde{\beta} - \beta) | \beta\} \\ &= E\left[\text{tr}\left\{A(\tilde{\beta} - \beta)(\tilde{\beta} - \beta)^T A^T\right\} | \beta\right] \\ &= \text{tr}\left[E\left\{A(\tilde{\beta} - \mu + \mu - \beta)(\tilde{\beta} - \mu + \mu - \beta)^T A^T | \beta\right\}\right]. \end{aligned}$$

- ▶ The expectation above is

$$A\left\{\text{var}(\tilde{\beta} | \beta) + (X^T X + V_*^-)^{-1} V_*^- (\beta - \beta_*)(\beta - \beta_*)^T V_*^- (X^T X + V_*^-)^{-1}\right\} A^T,$$

giving the MSE when estimating a fixed β .

- ▶ Taking expectations over the prior model for β gives

$$E\left(\|A\tilde{\beta} - A\beta\|^2\right) = \sigma^2 \text{tr}\left\{A(X^T X + V_*^-)^{-1} A^T\right\}, \quad (6)$$

which is larger than $A\text{var}(\tilde{\beta} | \beta)A^T$ and does not depend on β_* .

- ▶ This computation uses only the mean and variance, so holds under second-order assumptions.
- ▶ From now on we set $\beta_* = 0$, unless we state otherwise.

Equivalent degrees of freedom

- ▶ If we set $\beta_* = 0$, then the fitted values are

$$\tilde{y} = X\tilde{\beta} = X(X^T X + V_*^-)^{-1} X^T y = H_* y,$$

say.

- ▶ We define the **equivalent degrees of freedom** of the fit as

$$\mathbf{edf} = \text{tr}(H_*) = \text{tr}\{X(X^T X + V_*^-)^{-1} X^T\} = p - \text{tr}\{(X^T X + V_*^-)^{-1} V_*^-\},$$

- ▶ This is lower than p unless $V_*^- = 0$, so regularisation reduces the degrees of freedom by an amount that depends on V_* .
- ▶ The penalised estimate is a linear function of the unpenalised one (if it exists), as we can write

$$\tilde{\beta} = (X^T X + V_*^-)^{-1} X^T X \hat{\beta} = P_* \hat{\beta},$$

say. As

$$\mathbf{edf} = \text{tr}(H_*) = \text{tr}(P_*),$$

this gives an alternative formula useful in complex models.

How much penalisation?

- ▶ Often V_*^- depends on some $\lambda > 0$ that must be chosen, as well as σ^2 , which is usually estimated by a (penalised) residual sum of squares.
- ▶ To estimate λ , we compare y_j with its predicted value $\hat{y}_{\lambda,j}^- = x_j^T \hat{\beta}_{\lambda,-j}$, where $\hat{\beta}_{\lambda,-j}$ is

$$\hat{\beta}_{\lambda} = (X^T X + V_*^-)^{-1} X^T y$$

computed with the j th rows x_j and y_j of X and y omitted.

- ▶ Using Lemma 9, the **leave-one-out cross-validation** sum of squares is then

$$\text{CV}_{\lambda} = \sum_{j=1}^n (y_j - \hat{y}_{\lambda,j}^-)^2 = \|y - \hat{y}_{\lambda}^-\|^2 = \sum_{j=1}^n \frac{(y_j - \hat{y}_{\lambda,j}^-)^2}{(1 - h_{\lambda,j,j})^2},$$

where $\hat{y}_{\lambda,j}^-$ is the j th element of the complete-data fitted value $H_{\lambda}y$ and $h_{\lambda,j,j}$ is the j th diagonal element of $H_{\lambda} = X(X^T X + V_*^-)^{-1} X^T$ for the overall fit.

- ▶ More often we use the **generalized cross-validation** criterion

$$\text{GCV}_{\lambda} = \sum_{j=1}^n \frac{(y_j - \hat{y}_{\lambda,j}^-)^2}{\{1 - \text{tr}(H_{\lambda})/n\}^2}.$$

- ▶ Whichever criterion is used, it is typically minimised numerically over a grid of values of λ .

REML

- ▶ Cross-validation makes only second-order assumptions.
- ▶ Under normality, the marginal density of y is $\mathcal{N}\{X\beta_*, \sigma^2(I_n + XV_*X^T)\}$, so we could estimate β_* , σ^2 and λ by maximising the corresponding likelihood.
- ▶ If n and p are large, this results in biased estimates of λ and σ^2 , so we prefer to eliminate β_* , resulting in a **log restricted likelihood** whose form is given below, with $W_\lambda^{-1} = I_n + XV_*X^T$.

Lemma 20

In a model in which $y \sim \mathcal{N}(X\beta, \sigma^2 W_\lambda^{-1})$, where W_λ depends on a parameter λ , a log restricted likelihood for σ^2 and λ is

$$\ell_{\text{REML}}(\sigma^2, \lambda) \equiv \frac{1}{2} \log(|W_\lambda|/|X^T W_\lambda X|) - \frac{n-p}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y - \hat{y}_\lambda)^T W_\lambda (y - \hat{y}_\lambda),$$

where $\hat{\beta}_\lambda = (X^T W_\lambda X)^{-1} X^T W_\lambda y$ and $\hat{y}_\lambda = X \hat{\beta}_\lambda$.

For fixed λ the restricted maximum likelihood estimator of σ^2 is therefore

$$\hat{\sigma}_\lambda^2 = \frac{1}{n-p} (y - \hat{y}_\lambda)^T W_\lambda (y - \hat{y}_\lambda),$$

and the resulting profile log restricted likelihood for λ is

$$\ell_p(\lambda) \equiv \frac{1}{2} \log(|W_\lambda|/|X^T W_\lambda X|) - \frac{(n-p)}{2} \log \hat{\sigma}_\lambda^2.$$

Note on Lemma 20

- ▶ Suppose that $f(y; \alpha, \beta)$ depends on two parameters, that interest is focused on α , and that for fixed α there is a minimal sufficient statistic s_α for β . Then $f(y; \alpha, \beta) = f(y | s_\alpha; \alpha)f(s_\alpha; \alpha, \beta)$, and since the first density on the right is a proper conditional density not depending on β , we can use it for inference on α , in the form

$$\log f(y | s_\alpha; \alpha) = \log f(y; \alpha, \beta) - \log f(s_\alpha; \alpha, \beta).$$

As the left-hand side of this expression does not depend on β , we may be able to simplify the right-hand side by an astute choice of β .

- ▶ In the normal model we take $\alpha = (\sigma^2, \lambda)$. If α is fixed, then $s_\alpha = \hat{\beta}_\alpha = (X^T W_\lambda X)^{-1} X^T W_\lambda y$ is sufficient for β ; its distribution is $\mathcal{N}_p\{\beta, \sigma^2 (X^T W_\lambda X)^{-1}\}$. Hence

$$\ell_{\text{REML}}(\sigma^2, \lambda) = \log f(y \mid \hat{\beta}_\lambda; \sigma^2, \lambda) = \log f(y; \sigma^2, \lambda, \beta) - \log f(\hat{\beta}_\lambda; \sigma^2, \lambda, \beta)$$

which equals

$$\begin{aligned} -\frac{n}{2} \log \sigma^2 &+ \frac{1}{2} \log |W_\lambda| - \frac{1}{2\sigma^2} (y - X\beta)^T W_\lambda (y - X\beta) \\ &+ \frac{p}{2} \log \sigma^2 - \frac{1}{2} \log |X^T W_\lambda X| + \frac{1}{2\sigma^2} (\hat{\beta}_\lambda - \beta)^T X^T W_\lambda X (\hat{\beta}_\lambda - \beta), \end{aligned}$$

or equivalently, on setting $\beta = 0$ and $\hat{y}_\lambda = X\hat{\beta}_\lambda$,

$$\frac{1}{2} \log (|W_\lambda| / |X^T W_\lambda X|) - \frac{(n-p)}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y^T W_\lambda y - \hat{y}_\lambda^T X^T W_\lambda \hat{y}_\lambda).$$

- ▶ The last term reduces to the given form because $\hat{y}_\lambda^T W_\lambda (y - \hat{y}_\lambda) = 0$, so the term in brackets in the last displayed equation is the residual sum of squares $(y - \hat{y}_\lambda)^T W_\lambda (y - \hat{y}_\lambda)$.
- ▶ The restricted maximum likelihood estimator $\hat{\sigma}_\lambda^2$ and the profile log restricted likelihood for λ are obtained by maximising $\ell_{\text{REML}}(\sigma^2, \lambda)$, for fixed λ and then dropping constant terms from $\ell_{\text{REML}}(\hat{\sigma}_\lambda^2, \lambda)$.

Regularisation - Simple Applications

$\hat{\beta} \in \operatorname{argmin} \|y - X\beta\|^2 + P_\lambda(\beta)$ where $P_\lambda(\beta)$ is the penalty function depending on $\lambda > 0$ and

- $P_\lambda(\beta) = \lambda \|\beta\|_2$ \leadsto ridge regression
- $= \lambda \|\beta\|_1$ \leadsto Lasso regression
- $= \alpha \|\beta\|_2 + \gamma \|\beta\|_1$ \leadsto Elastic net.

Last session: Motivate the use of penalization with Bayesian calculus:

We showed that:

$$\begin{cases} y | \beta, \sigma^2 \sim \mathcal{N}(X\beta, \sigma^2 I) \\ \beta | \sigma^2 \sim \mathcal{N}(\beta_0, \sigma^2 V_0) \end{cases}$$

then the posterior distribution of $\beta | y, \sigma^2$ is also Gaussian

$$\tilde{y} = X\tilde{\beta} = \underbrace{X^T(X^T X + V_0)^{-1}}_{= H_0} (X^T X \hat{\beta} + V_0^{-1} \beta_0)$$
 if $\beta_0 = 0$ then \tilde{y} is just a linear transformation of \hat{y}
 (OLS: $\hat{y} = X\hat{\beta}$)

Equivalent degrees of freedom

$$\begin{cases} \tilde{y} = H_0 \hat{y} \\ \text{edf} = \operatorname{tr}(H_0) = p - \operatorname{tr}((X^T X + V_0)^{-1} V_0) \end{cases}$$

\hookrightarrow Define CV, GCV

Penalization is used when # covariate parameters \gg # observations

$$\hat{\beta} = (X^T X)^{-1} X^T y$$
 \hookrightarrow ill defined.

Ridge regression

- ▶ Used for prediction when X is close to singular.
- ▶ If the first column of X is 1_n , we set $\beta_* = 0$ and $V_*^- = \lambda S = \lambda \text{diag}(0, I_{p-1})$, giving

$$\hat{\beta}_\lambda = (X^T X + \lambda S)^{-1} X^T y, \quad \hat{y}_\lambda = X \hat{\beta}_\lambda = \underbrace{X(X^T X + \lambda S)^{-1} X^T}_{H_\lambda} y = H_\lambda y,$$

and effective degrees of freedom

$$\mathbf{edf}_\lambda = \text{tr}(H_\lambda) = \text{tr}\{(X^T X + \lambda S)^{-1} X^T X\} = \sum_{r=1}^p \frac{1}{1 + \lambda \delta_r}, = 1 + \sum_{r=2}^p \frac{1}{1 + \lambda \delta_r}$$

where $\delta_p \geq \dots \geq \delta_2 > \delta_1 = 0$ are the eigenvalues of $(X^T X)^{-1/2} S (X^T X)^{-1/2}$.

- ▶ As λ increases from zero to infinity, \mathbf{edf}_λ decreases from $p = \mathbf{rank}(X)$ to 1. The two are equivalent, but \mathbf{edf}_λ is more easily interpreted, because it is not related to the scale of X .
- ▶ The inverse exists even if $X^T X$ is singular, but if it is invertible then

$$\hat{\beta}_\lambda = (X^T X + \lambda S)^{-1} (X^T X + \lambda S - \lambda S) (X^T X)^{-1} X^T y = \hat{\beta} - \lambda (X^T X + \lambda S)^{-1} S \hat{\beta},$$

so as $\lambda \rightarrow \infty$ all the elements of $\hat{\beta}_\lambda$ tend to zero, other than the first. This corresponds to reducing the prior variance to zero, thereby giving the data themselves less and less influence on the elements of $\hat{\beta}_\lambda$ other than the first, and thus stabilises the estimator.

In practice

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

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

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

- ▶ Write $X = (1_n \ X')$, $\beta = (\beta_0 \ \gamma)^\top$
- ▶ Recentre/rescale the covariates (columns) defining:
$$Z_j = \frac{\sqrt{n}}{\text{sd}(X'_j)} (X'_j - 1_n \overline{X'_j})$$
 - ▶ Coefficients now have common scale
 - ▶ Interpretation of β_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 Z_j are all orthogonal to 1_n and are of unit norm.

- ▶ Since $Z_j \perp 1_n$ for all, j , we can estimate β_0 and γ by two separate regressions (orthogonality).
- ▶ Least squares estimators based on the *standardised design matrix* become

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

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

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

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

→ Motivation of adding the ridge term λI

↔ Can see that $(\hat{\beta}_0 \quad \hat{\gamma}) = (\bar{y} \quad (Z^T Z + \lambda I)^{-1} Z^T y)$ minimizes

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

or equivalently

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

instead of least squares estimator which minimizes

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

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!

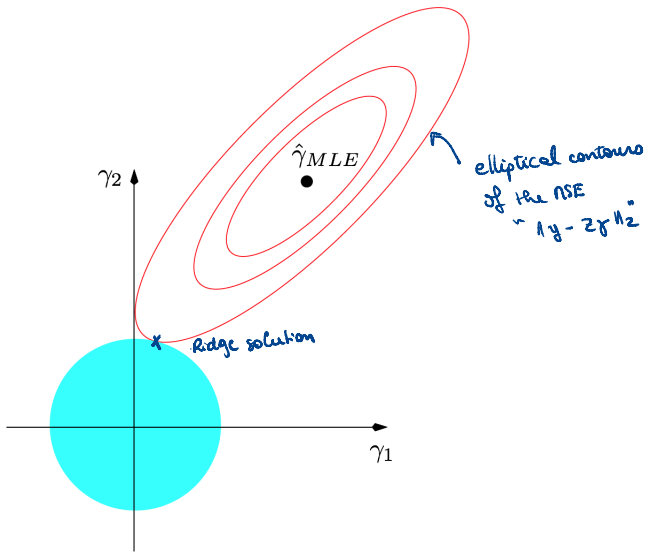


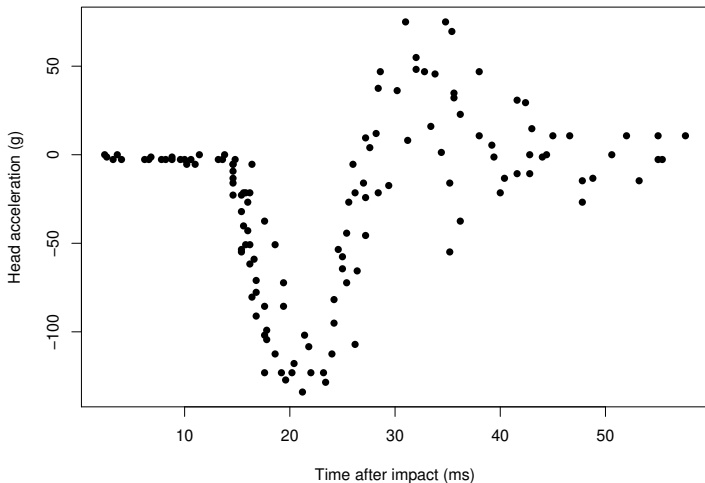
Figure: L^2 Shrinkage [Ridge Regression]

Semiparametric regression

- ▶ Normal linear model has two main aspects:
 - ▶ **systematic variation**, $E(y) = \mu$, and $\mu = X\beta$ with parameters β ;
 - ▶ **stochastic variation**, $y \sim \mathcal{N}_n(\mu, \sigma^2 I_n)$.
- ▶ Can relax the stochastic assumption using other distributions or second-order assumptions, but still have parametric model for the systematic part.
- ▶ Often want to relax systematic part for more flexible models, for
 - ▶ exploratory data analysis — ‘will a linear model be adequate?’
 - ▶ confirmatory data analysis — ‘I’ve fitted a linear model, is it adequate?’
 - ▶ general modelling — ‘the data are too complex to expect a simple parametric model to work, so what can I do?’
 - ▶ semiparametric modelling — ‘I will use a parametric model for the effects of interest, but can I model nuisance effects more flexibly?’
- ▶ Most basic tool is the **scatterplot smoother**.

Example: Motorcycle data

Measurements of head acceleration (g) at time after impact (ms) in a simulated motorcycle accident, used to test crash helmets:



Scatterplot smoothing

- ▶ Have data $(x_1, y_1), \dots, (x_n, y_n)$, with $x_- \leq x_1 < \dots < x_n \leq x_+$ (ahem) and we wish to estimate $E(y) = \mu(x)$, for $x \in \mathcal{X} = [x_-, x_+]$.
- ▶ Suppose that $\mu \in \mathcal{M}$, a function space spanned by n linearly independent basis functions that can be identified by evaluation at x_1, \dots, x_n , and let $\mu_j = \mu(x_j)$.
- ▶ Can choose a basis $\{b_1(x), \dots, b_n(x)\}$ for \mathcal{M} such that $\mu(x) = \sum_{j=1}^n \mu_j b_j(x)$ interpolates $(x_1, \mu_1), \dots, (x_n, \mu_n)$.
- ▶ Suppose that \mathcal{M} contains the linear functions on \mathcal{X} and that the second derivatives of the $b_j(x)$ are not all zero, so functions in \mathcal{M} may also be nonlinear in x .
- ▶ To estimate μ we minimise a **penalised sum of squares**,

$$\sum_{j=1}^n \{y_j - \mu(x_j)\}^2 + \lambda \int_{\mathcal{X}} \{\mu''(x)\}^2 dx, \quad (7)$$

where the **roughness penalty** imposes smoothness: if $\lambda \rightarrow 0$, then $\mu(x_j) \rightarrow y_j$ and $\hat{\mu}$ interpolates, but when $\lambda \rightarrow \infty$ even tiny wiggles in μ will give a huge penalty, making $\hat{\mu}$ linear.

- ▶ The penalty does not affect linear functions, so $\mathcal{M} = \mathcal{L} \oplus \mathcal{P}$, where \mathcal{L} and \mathcal{P} are the two-dimensional vector space of linear functions on \mathcal{X} and an $(n - 2)$ -dimensional vector space of nonlinear functions on \mathcal{X} , and \oplus denotes addition of vector spaces.

- ▶ The roughness term is

$$\int_{\mathcal{X}} \{\mu''(x)\}^2 dx = \int_{\mathcal{X}} \left\{ \sum_{j=1}^n \mu_j b_j''(x) \right\}^2 dx = \sum_{i,j=1}^n \mu_i \mu_j \int_{\mathcal{X}} b_i''(x) b_j''(x) dx = \mu^T S \mu,$$

say, where $\mu^T = (\mu_1, \dots, \mu_n)$.

- ▶ $S_{n \times n}$ has (i, j) element $\int_{\mathcal{X}} b_i''(x) b_j''(x) dx$ and is symmetric and positive semi-definite of rank $n - 2$, because linear functions are unpenalised, so $S1_n = S(x_1, \dots, x_n)^T = 0$.
- ▶ The penalised sum of squares

$$(y - \mu)^T (y - \mu) + \lambda \mu^T S \mu \equiv -2\mu^T y + \mu^T (I_n + \lambda S) \mu,$$

is minimised by $\hat{\mu}_\lambda = (I_n + \lambda S)^{-1} y$.

- ▶ As λ increases from zero, the fitted value $\hat{\mu}_\lambda$ shrinks from y towards the straight-line regression fit to y , which is unpenalised.
- ▶ The equivalent degrees of freedom are $\mathbf{edf}_\lambda = \text{tr}(H_\lambda) = \sum_{j=1}^n (1 + \lambda \delta_j)^{-1}$, where $\delta_1 \geq \dots \geq \delta_3 > \delta_2 = \delta_1 = 0$ are the eigenvalues of S . As λ increases \mathbf{edf}_λ decreases monotonically from $\mathbf{edf}_0 = n$ towards $\mathbf{edf}_\infty = 2$.

- ▶ In principle we might take any basis functions, but in practice we usually take local polynomials known as **splines** that have good approximation properties.
- ▶ There are many forms of splines, which
 - ▶ are often cubic polynomials with finite support between values of x known as **knots**, x_1^*, \dots, x_K^* , and then S is tri-diagonal,
 - ▶ sometimes form a **natural cubic spline**, which has $K = n$ and certain optimality properties,
 - ▶ are discussed in more detail later.
- ▶ If there is no penalisation ($\lambda = 0$) then we have a standard linear model, and spline basis functions are called **regression splines**.
- ▶ Under second-order assumptions we choose λ by minimising $\text{CV}(\lambda)$ or $\text{GCV}(\lambda)$.
- ▶ Under normal-theory assumptions we can use REML to estimate σ^2 and λ .
- ▶ Obvious generalisation allows weight matrix $W = \text{diag}(w_1, \dots, w_n)$.
- ▶ If the x_1, \dots, x_n are not unique, write $E(y) = N_{n \times n'} \mu_{n' \times 1}$ in terms of the means μ at the n' unique elements of x , and minimise

$$(y - N\mu)^T W (y - N\mu) + \lambda \mu^T S \mu.$$

where $S_{n' \times n'}$ arises as before from the roughness penalty on $\mu(x)$.

Linear, quadratic and cubic B -splines

