

# Regularized estimation of LM and GLM

## Part 1: Ridge regression

Pedro Delicado

Departament d'Estadística i Investigació Operativa  
Universitat Politècnica de Catalunya

## 1 Introduction

## Multiple linear regression model

## 2 Ridge regression

## Linear estimators of a regression function

## Choosing the tuning parameter

## Cross-Validation and Generalized Cross-Validation

## REFERENCES:

Section 3.4 in Hastie, Tibshirani, and Friedman (2009)

## 1 Introduction

## Multiple linear regression model

## Linear estimators of a regression function

## Choosing the tuning parameter

## Cross-Validation and Generalized Cross-Validation

# Introduction

- In the multiple linear regression model (with  $n$  observations and  $p$  predictors,  $p$  possibly greater than  $n$ ) we consider the penalized least squares coefficients estimator where the penalization is given by the  $L_1$  or the  $L_2$  norms of the estimator.
- This procedure leads to **ridge regression** ( $L_2$  penalization) and to **lasso** (least absolute shrinkage and selection operator) estimation ( $L_1$  penalization).
- In the pathway, we will learn:
  - Linear estimators of a regression function.
  - Effective number of parameters (or effective degrees of freedom) of a regression estimator.
  - Tuning parameters choice based on leave-one-out cross-validation,  $k$ -fold cross-validation or generalized cross-validation.
  - Efficient computation of leave-one-out cross-validation for linear estimators.

## 1 Introduction

## Multiple linear regression model

## Linear estimators of a regression function

## Choosing the tuning parameter

## Cross-Validation and Generalized Cross-Validation

## Multiple linear regression model

- Consider that  $n$  pairs  $(x_i, y_i)$ ,  $i = 1, \dots, n$  of data,  $y_i \in \mathbb{R}$  and  $x_i \in \mathbb{R}^p$ , are observed from the **multiple linear regression model**

$$y_i = \beta_0 + \sum_{j=1}^p x_{ij}\beta_j + \varepsilon_i,$$

where  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d. r.v. with zero mean and variance  $\sigma^2$ , and  $\beta = (\beta_0, \dots, \beta_p)^\top \in \mathbb{R}^{p+1}$  is a vector of unknown coefficients.

- Fitting the model consists in providing estimators for  $\beta$  and  $\sigma^2$ , jointly with information about the sampling distribution of these estimators (standard errors, hypothesis testing, among others).

## Ordinary Least Squares (OLS)

- Ordinary Least Squares (OLS) estimator:

$$\hat{\beta}_{\text{OLS}} = \arg \min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

- In matrix notation:  $\hat{\beta}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ .
- $\hat{\beta}_{\text{OLS}}$  is an unbiased estimator of  $\beta$ .



- **Gauss–Markov Theorem.** For any  $\mathbf{a} \in \mathbb{R}^{p+1}$ , the OLS estimator of the linear combination  $\mathbf{a}^\top \boldsymbol{\beta}$ , namely  $\mathbf{a}^\top \hat{\boldsymbol{\beta}}_{\text{OLS}}$ , is unbiased and it has the lowest variance among the linear unbiased estimates of  $\mathbf{a}^\top \boldsymbol{\beta}$ .
- In particular, following the Bayes rule, the prediction for a new observation  $\mathbf{x}$ , is  $\hat{y} = \mathbf{x}^\top \boldsymbol{\beta}$ .
- So its best unbiased estimator is  $\hat{y}_{\text{OLS}} = \mathbf{x}^\top \hat{\boldsymbol{\beta}}_{\text{OLS}}$ .

<sup>1</sup>Slides with colored background are optional material.

## Multicollinearity and bad conditioned matrices

- The computation of  $\hat{\beta}_{OLS}$  is numerically unstable when  $\mathbf{X}^T \mathbf{X}$  is close to be singular.
- **Condition number** of a symmetric matrix  $\mathbf{A}$ :  $\kappa(\mathbf{A}) = \sqrt{\frac{\gamma_{\max}}{\gamma_{\min}}}$ , where  $\gamma_{\max}$  and  $\gamma_{\min}$  are, respectively, the largest and lowest eigenvalue absolute values of  $\mathbf{A}$ .
- $\mathbf{A}$  is not invertible if and only if  $\kappa(\mathbf{A}) = \infty$ .
- A large value of  $\kappa(\mathbf{A})$  (in practice, larger than 30), indicates that numerical problems may appear when inverting  $\mathbf{A}$ .
- In these cases we say that  **$\mathbf{A}$  is bad conditioned**.
- If  $\mathbf{X}^T \mathbf{X}$  is bad conditioned: unstable computation of  $\hat{\beta}_{OLS}$ .
- A large condition number indicates that  $\mathbf{X}$  is close to be singular, that is, close that some columns of  $\mathbf{X}$  can be written as linear combinations of the other.
- We talk about **multicollinearity** between columns of  $\mathbf{X}$ .

## Regularized regression

- Beyond numerical problems,  $\hat{\beta}_{OLS}$  can not be computed when the rank of  $\mathbf{X}$  is lower than the number of variables ( $p + 1$ ) (this is an extreme case of multicollinearity).
- This is the case when  $p + 1 > n$  (or  $p + 1 \gg n$ , as it can happen in applications with large scale data).
- In practical terms, what happens is that  $\mathbf{y}$  can be written as a linear combination of the predictors using infinitely many coefficient vectors, for which the objective OLS objective function is equal to 0, its minimum. So there is no way to select *the best* among those coefficient vectors.
- **Shrinkage (or regularized) methods:** They add a penalty (depending on  $\beta$ ) to the objective function in such a way that the new optimum is attained at a unique vector  $\hat{\beta}$ .
- The unbiasedness of OLS estimation is lost, but the new estimators may have lower Mean Square Error (and they are numerically stable).

## Multiple linear regression model

## 2 Ridge regression

## Linear estimators of a regression function

## Choosing the tuning parameter

## Cross-Validation and Generalized Cross-Validation

## Ridge regression

- The ridge coefficients minimize a penalized sum of squares residuals:

$$\begin{aligned}\hat{\beta}_{\text{ridge}} &= \arg \min_{\beta} \text{SSR}_{\text{pen}}(\beta) \\ &= \arg \min_{\beta} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\} \\ &= \arg \min_{\beta} \{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \|\beta_{-0}\|_2^2 \}\end{aligned}$$

where  $\beta_{-0} = (\beta_1, \dots, \beta_p)^\top$ .

- Here  $\lambda \geq 0$  is a complexity parameter that controls the amount of shrinkage: the larger the value of  $\lambda$ , the greater the amount of shrinkage of  $\beta_{-0}$  toward zero.

- Ridge regression is a penalized least squares problem:

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

- Alternative expression, constrained least squares problem:

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 \leq t,$$

for  $t \geq 0$ . There is a one-to-one decreasing correspondence between parameters  $\lambda \in [0, \infty)$  and  $t \in (0, \|\hat{\beta}_{-0, \text{OLS}}\|^2]$ .

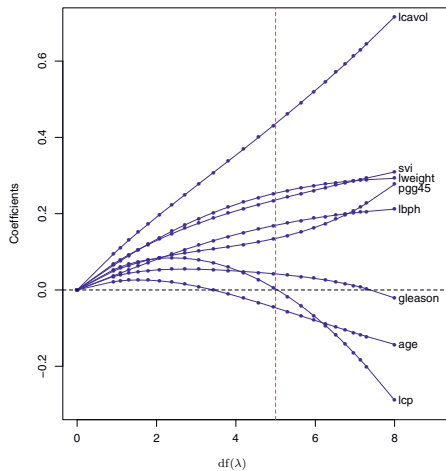
- Observe that changes in scale of the explanatory variables affect the constraint effects (or, equivalently, the effects of penalization term).
- For this reason, from now on we assume that the predictor variables have zero mean and unit variance (otherwise we center and standardize them in advance):

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

- Moreover, the response variable is assumed to have zero mean ( $\sum_{i=1}^n y_i = 0$ ), that is,  $\beta_0 = 0$ .







Source: Hastie, Tibshirani, and Friedman (2009)

## Explicit solution for the ridge regression

- The ridge regression estimators are the solution of the penalized least squares problem

$$\min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n \left( y_i - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

that can be expressed as

$$\min_{\beta \in \mathbb{R}^p} \Psi(\beta) = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^\top \beta,$$

that has an explicit solution, as we show now.

- Taking the gradient

$$\nabla \Psi(\beta) = -2\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\beta) + 2\lambda\beta,$$

and solving in  $\beta$  the equation  $\nabla\Psi(\beta) = \mathbf{0}$ , we obtain

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X} + \lambda I_p)^{-1} \mathbf{X}^\top \mathbf{y}.$$

- Ridge regression estimator:  $\hat{\beta}_{\text{ridge}} = \left( \mathbf{X}^T \mathbf{X} + \lambda I_p \right)^{-1} \mathbf{X}^T \mathbf{y}$ .
- Therefore, for any  $\mathbf{x} \in \mathbb{R}^p$ , the corresponding predicted value is

$$\hat{y} = \mathbf{x}^\top \hat{\boldsymbol{\beta}}_{\text{ridge}} = \mathbf{x}^\top \left( \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p \right)^{-1} \mathbf{X}^\top \mathbf{y}.$$

- The vector of fitted values is

$$\hat{\mathbf{y}} = \mathbf{X} \left( \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p \right)^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{H}_\lambda \mathbf{y}.$$

- Compare with the OLS solution:  $\hat{\beta}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ .

$$\hat{\mathbf{y}}_{\text{OLS}} = \mathbf{X} \left( \mathbf{X}^{\top} \mathbf{X} \right)^{-1} \mathbf{X}^{\top} \mathbf{y} = \mathbf{H} \mathbf{y},$$

where  $H = X(X^T X)^{-1} X^T$  is called the **hat matrix**.

- $\lim_{\lambda \rightarrow 0} \hat{\beta}_{\text{ridge}} = \hat{\beta}_{\text{OLS}}, \lim_{\lambda \rightarrow \infty} \hat{\beta}_{\text{ridge}} = \mathbf{0}.$





As we are assuming that the explanatory variables have zero mean, we have that  $\mathbf{X}^T \mathbf{X}$  is the sample covariance matrix. Then the columns of  $\mathbf{V}$  are the **principal components** of  $\mathbf{X}$ . Moreover the columns of  $\mathbf{UD}$  are the scores of the observed data in the principal components.

## Numerical stability of ridge regression

- $\hat{\beta}_{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$
- Let us compute the condition number of  $\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p$ ,

$$\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p = \mathbf{V} \mathbf{D}^2 \mathbf{V}^T + \lambda \mathbf{V} \mathbf{V}^T = \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}_p) \mathbf{V}^T.$$

- $(\mathbf{D}^2 + \lambda \mathbf{I}_p)$  is a diagonal matrix whose elements in the diagonal are

$$d_j^2 + \lambda = \gamma_j + \lambda, \quad j = 1, \dots, p.$$

- Therefore the condition number of  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p$  is

$$\kappa\left(\mathbf{X}^{\top}\mathbf{X}+\lambda\mathbf{I}_p\right)=\sqrt{\frac{\gamma_1+\lambda}{\gamma_p+\lambda}}$$

lower than  $\kappa(\mathbf{X}^\top \mathbf{X}) = \sqrt{\gamma_1/\gamma_p}$  for all  $\lambda > 0$ .

- By the way,  $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} = \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} \mathbf{V}^T$ , and  $(\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} = \text{Diagonal}(1/(d_j^2 + \lambda), j = 1, \dots, p)$ .

## Variance of the ridge regression estimator

Remember that  $\hat{\beta}_{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$ . Then

$$\begin{aligned}\text{Var}(\hat{\beta}_{\text{ridge}}) &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \text{Var}(\mathbf{y}) \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1}.\end{aligned}$$

From the s.v.d. of  $\mathbf{X}$ ,  $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ , we have deduced that  $\mathbf{X}^\top\mathbf{X} = \mathbf{V}\mathbf{D}^2\mathbf{V}^\top$  and that  $(\mathbf{X}^\top\mathbf{X} + \lambda\mathbf{I}_p)^{-1} = \mathbf{V}(\mathbf{D}^2 + \lambda\mathbf{I}_p)^{-1}\mathbf{V}^\top$ . Therefore,

$$\begin{aligned}\text{Var}(\hat{\beta}_{\text{ridge}}) &= \sigma^2 \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} \mathbf{D}^2 (\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} \mathbf{V}^\top \\ &= \sigma^2 \mathbf{V} \text{Diagonal} (d_j^2 / (d_j^2 + \lambda)^2, j = 1, \dots, p) \mathbf{V}^\top.\end{aligned}$$



## Multiple linear regression model

## 2 Ridge regression

## Linear estimators of a regression function

## Choosing the tuning parameter

## Cross-Validation and Generalized Cross-Validation

## Linear estimators of a regression function

- Let  $(x_i, y_i)$ ,  $i = 1, \dots, n$ , be  $n$  i.i.d. observs. from the r.v.  $(\mathbf{X}, Y)$ .
- Let  $m(x) = \mathbb{E}(Y|\mathbf{X} = x)$  be the regression function of  $Y$  over  $\mathbf{X}$ .
- Let  $\hat{m}(x)$  an estimator (parametric, non-parametric, ...) of the regression function  $m(x)$ .
- We say that  $\hat{m}(x)$  is a **linear estimator** when for any fix  $x$ ,  $\hat{m}(x)$  is a linear function of  $y_1, \dots, y_n$ :

$$\hat{m}(x) = \sum_{i=1}^n w_i(x) y_i,$$

where in fact  $w_i(\mathbf{x}) = w_i(\mathbf{x}; \mathbf{x}_1, \dots, \mathbf{x}_n)$ .

- For the  $n$  observed values  $x_i$  of the explanatory variable, let

$$\hat{y}_i = \hat{m}(x_i) = \sum_{j=1}^n w_j(x_i) y_j$$

be the fitted values.

- In matrix format,

$$\hat{y} = Wy,$$

where the column vectors  $\mathbf{y}$  and  $\hat{\mathbf{y}}$  have elements  $y_i$  and  $\hat{y}_i$ , respectively, and the matrix  $\mathbf{W}$  has generic  $(i, j)$  element

$$w_{ij} = w_j(x_i).$$

- The matrix  $\mathbf{W}$  is analogous to the **hat matrix**  $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  in OLS estimation of the multiple linear regression:

$$\hat{y}_{OLS} = X(X^T X)^{-1} X^T y = Hy.$$

- Observe that ridge regression is a linear estimation method:

$$\hat{\mathbf{y}}_{\text{ridge}} = \mathbf{X} \left( \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p \right)^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}_\lambda \mathbf{y}.$$

## Effective number of parameters for linear estimators

- Consider the multiple linear regression with  $p$  regressors (including the constant term, if it appears in the model):

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \varepsilon,$$

$\mathbf{X}$  being a  $n \times p$  matrix,  $\boldsymbol{\beta} \in \mathbb{R}^p$ .

- It is known that

$$\text{Trace}(\mathbf{H}) = \text{Trace}(\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) = \text{Trace}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}) = \text{Trace}(\mathbf{I}_p) = p,$$

that is the number of parameters in the model.

- For a linear estimator with matrix  $\mathbf{W}$  ( $\hat{\mathbf{y}} = \mathbf{W}\mathbf{y}$ ) we define

$$\nu = \text{Trace}(\mathbf{W}) = \sum_{i=1}^n w_{ii},$$

the sum of diagonal elements of  $\mathbf{W}$ .

- $\nu = \text{Trace}(\mathbf{W})$  is called the **effective number of parameters** of the linear estimator corresponding to matrix  $\mathbf{W}$ .
- In some books (and softwares)  $\nu$  is called **effective degrees of freedom** (df) of the regression estimator. This is the terminology used by Hastie, Tibshirani, and Friedman (2009) and Hastie, Tibshirani, and Wainwright (2015), and related packages.
- The interpretation of  $\nu$  as the effective number of parameters is valid for any linear estimator of the regression function (parametric, nonparametric, ...).
- Then we can compare the degree of complexity of two linear estimators of a regression function just comparing their effective numbers of parameters.
- Usually a good estimator of  $\sigma^2$ , the residual variance, is

$$\hat{\sigma}^2 = \frac{1}{n - \nu} \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

## Effective number of parameters in ridge regression

In the case of ridge regression  $\nu = \nu(\lambda) = \text{df}(\lambda)$  has an explicit expression:

$$\mathbf{W} = \mathbf{H}_\lambda = \mathbf{X} \left( \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p \right)^{-1} \mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top \mathbf{V} \left( \mathbf{D}^2 + \lambda \mathbf{I}_p \right)^{-1} \mathbf{V}^\top \mathbf{V} \mathbf{D} \mathbf{U}^\top =$$

$$\mathbf{U} \mathbf{D} \left( \mathbf{D}^2 + \lambda \mathbf{I}_p \right)^{-1} \mathbf{D} \mathbf{U}^\top = \mathbf{U} \left( \text{Diagonal}(d_j^2 / (d_j^2 + \lambda), j = 1, \dots, p) \right) \mathbf{U}^\top$$

$$\Rightarrow \nu(\lambda) = \text{df}(\lambda) = \text{Trace}(\mathbf{H}_\lambda) =$$

$$\text{trace}(\mathbf{U} \left( \text{Diagonal}(d_j^2 / (d_j^2 + \lambda), j = 1, \dots, p) \right) \mathbf{U}^\top) =$$

$$\text{trace}(\left( \text{Diagonal}(d_j^2 / (d_j^2 + \lambda), j = 1, \dots, p) \right) \mathbf{U}^\top \mathbf{U}) =$$

$$\text{trace}(\left( \text{Diagonal}(d_j^2 / (d_j^2 + \lambda), j = 1, \dots, p) \right)) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}.$$

$$\nu(\lambda) = \text{df}(\lambda) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$$

- $\lim_{\lambda \rightarrow \infty} \text{df}(\lambda) = 0$ ,  $\lim_{\lambda \rightarrow 0} \text{df}(\lambda) = \text{rank}(\mathbf{X})$ .
- The effective number of parameters  $\nu(\lambda) = \text{df}(\lambda)$  is a decreasing function of penalizing parameter  $\lambda$ :
  - Small values of  $\lambda$  correspond to large numbers  $\nu$  of effective parameters, close to the number of linearly independent explanatory variables (usually  $\min\{n, p\}$ ), allowing complex and flexible estimators.
  - Large values of  $\lambda$  correspond to small numbers  $\nu$  of effective parameters, that is, to regression estimators with low complexity and flexibility.

## Practice:

- Prostate data: Effective number of parameters in ridge regression.
- Use the R script `prostate.ridge.regression.R`.



## Effective degrees of freedom for non-linear estimators

- Let  $\hat{m}(x)$  an estimator of the regression function  $m(x)$  (a random function because it is based on  $(Y_1, \dots, Y_n)$ ). Let  $\hat{Y}_i = \hat{m}(x_i)$ .
- The effective degrees of freedom of  $\hat{m}(x)$  is defined as

$$\text{df}(\hat{m}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(\hat{Y}_i, Y_i).$$

- Interpretation:
  - A very flexible regression estimator  $\hat{m}(x)$  will be able to interpolate the observed data, and then
 
$$\hat{Y}_i = Y_i, \text{Var}(\hat{Y}_i) = \text{Var}(Y_i) = \sigma^2, \text{Cov}(\hat{Y}_i, Y_i)/\sigma^2 = \text{Cor}(\hat{Y}_i, Y_i) = 1,$$
 so  $\text{df}(\hat{m}) = n$ :  $\hat{m}(x)$  has as many degrees of freedom as the number of observed data.
  - The constant function equal to the sample mean of  $Y_1, \dots, Y_n$  for all  $x$  has 1 degree of freedom.
  - A function that is constant in  $x$  has 0 degrees of freedom if this constant does not depend on the data.

## Both definitions of df coincide in linear estimators

Assume that  $\hat{m}(x)$  is a linear estimator with matrix  $\mathbf{W}$ . Assume also that  $\mathbb{E}(\mathbf{Y}) = \mathbf{0}$ . Then

$$\begin{aligned} \text{df}(\hat{m}) &= \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(\hat{Y}_i, Y_i) = \frac{1}{\sigma^2} \text{Trace} \left( \text{Cov}(\hat{\mathbf{Y}}, \mathbf{Y}) \right) = \\ &= \frac{1}{\sigma^2} \text{Trace} \left( \mathbb{E}(\hat{\mathbf{Y}} \mathbf{Y}^\top) \right) = \frac{1}{\sigma^2} \text{Trace} \left( \mathbb{E}(\mathbf{W} \mathbf{Y} \mathbf{Y}^\top) \right) = \\ &= \frac{1}{\sigma^2} \text{Trace} \left( \mathbf{W} \mathbb{E}(\mathbf{Y} \mathbf{Y}^\top) \right) = \frac{1}{\sigma^2} \text{Trace} (\mathbf{W} \sigma^2 \mathbf{I}_p) = \text{Trace}(\mathbf{W}). \end{aligned}$$

## 1 Introduction

Multiple linear regression model

## 2 Ridge regression

Linear estimators of a regression function

**Choosing the tuning parameter**

Cross-Validation and Generalized Cross-Validation

## Choosing the tuning parameter $\lambda$

- The tuning parameter  $\lambda$  can be chosen by cross-validation (CV),  $k$ -fold cross-validation ( $k$ -fold CV) or by generalized cross-validation (GCV).
- Given the expression of  $\hat{\beta}_{\text{ridge}}$  (linear in  $\mathbf{y}$ ) CV and GCV are not computationally expensive.
- We will first introduce these concepts before talking about efficient computation.

- **Predictive Mean Square Error (PMSE)**. It is the expected squared error made when predicting

$$Y = m(x) + \varepsilon$$

by  $\hat{m}(x)$ , where  $x$  is an observation of the random variable  $\mathbf{X}$ , distributed as the observed explanatory variable, when  $\mathbf{X}$  and  $\varepsilon$  are independent from the sample  $\mathcal{Z} = \{(\mathbf{X}_i, Y_i) : i = 1, \dots, n\}$  used to compute  $\hat{m}$ :

$$\text{PMSE}(\hat{m}) = \mathbb{E}_{\mathcal{Z}, \mathbf{X}, \varepsilon} [(Y - \hat{m}(\mathbf{X}))^2].$$

## Prediction error in a validation set

- When the number of available data is large (as it usually happens in data mining or in Big Data problems) the sample is randomly divided in three sets:
  - The **training set**: it is used to fit the model.
  - The **validation set**: it is used to compute feasible versions of the PMSE for model selection and/or parameter tuning.
  - The **test set**: it is used to evaluate the generalization (or prediction) error of the final chosen model in independent data.
- Assuming that at least a validation set has been preserved, an estimation of PMSE is the **Predictive Mean Squared Error in the validation set**:

$$\text{PMSE}_{\text{val}}(\hat{m}) = \frac{1}{n_V} \sum_{i=1}^{n_V} (y_i^V - \hat{m}(x_i^V))^2,$$

where  $(x_i^V, y_i^V)$ ,  $i = 1, \dots, n_V$ , is the validation set and  $\hat{m}(x)$  is the estimator computed using the training set.

## Multiple linear regression model

## 2 Ridge regression

## Linear estimators of a regression function

## Choosing the tuning parameter

## Cross-Validation and Generalized Cross-Validation

## Leave-one-out cross-validation

- When the sample size does not allow us to set a validation set aside, **leave-one-out cross-validation** is an attractive alternative:

- 1 Remove the observation  $(x_i, y_i)$  from the sample and fit the regression using the other  $(n - 1)$  data. Let  $\hat{m}_{(i)}(x)$  be the resulting estimator.
- 2 Now use  $\hat{m}_{(i)}(x_i)$  to predict  $y_i$ .
- 3 Repeat the previous steps for  $i = 1, \dots, n$ .
- 4 Compute

$$\text{PMSE}_{\text{CV}}(\hat{m}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{m}_{(i)}(x_i))^2.$$

- In ridge regression:

$$\lambda_{\text{CV}} = \arg \min_{\lambda \geq 0} \text{PMSE}_{\text{CV}}(\lambda) = \frac{1}{n} \sum_{i=1}^n (y_i - x_i^{\top} \hat{\beta}_{\text{ridge}, \lambda}^{(i)})^2.$$



## Practice:

- Prostate data: Leave-one-out cross-validation in ridge regression.
- Use the R script `prostate.ridge.regression.R`.

## $k$ -fold cross validation

- $\text{PMSE}_{\text{CV}}(\hat{m})$  is an approximately unbiased estimator of  $\text{PMSE}(\hat{m})$ , but has a considerable variance.
- The variance can be reduced doing  $k$ -fold cross-validation: The sample is randomly divided in  $k$  subsets, each of them is removed by turns from the sample, the model is estimated with the other  $(k - 1)$  subsamples and the removed subsample is used to compute prediction errors.
- $n$ -fold cross-validation is leave-one-out cross-validation.
- $k$ -fold cross-validation has lower variance than leave-one-out cross-validation but larger bias.
- General recommendation: Use 5-fold or 10-fold cross-validation.

## Efficient computation of $\text{PMSE}_{\text{CV}}$

- Consider a **linear estimator** of the regression function with matrix  $W = (w_{ij})_{i,j}$ :  $\hat{y} = Wy$ .
- That is

$$\hat{y}_i = \sum_{j=1}^n w_{ij} y_j, \quad i = 1, \dots, n,$$

where  $w_{ij} = w_j(\mathbf{x}_i) = w_j(\mathbf{x}_i; \mathbf{x}_1, \dots, \mathbf{x}_n)$ .

- In these cases  $\text{PMSE}_{\text{CV}}$  can be calculated avoiding the computational cost of fitting  $n$  different regression models.
- For most linear estimators it can be proved that

$$\text{PMSE}_{\text{CV}} = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{y}_i}{1 - w_{ii}} \right)^2.$$

## Self-stable linear estimators

(See Fan, Li, Zhang, and Zou 2020, page 48)

- Let  $\hat{m}(x)$  be a linear estimator of the regression function  $m(x) = \mathbb{E}(Y|X = x)$  fitted on the data  $(x_i, y_i)$ ,  $i = 1, \dots, n$ .
- Let  $x_0$  be a new covariate vector and  $\hat{m}(x_0)$  be its predicted value using the linear estimator.
- We augment the data set by including  $(x_0, \hat{m}(x_0))$  as a new point, and refit the linear estimator on this augmented data set:  $\hat{m}_a(x)$ .
- The linear estimator of the regression function is said to be **self-stable** if the fit based on the augmented data set is identical to the fit based on the original data regardless of  $x_0$ .

## Self-stability and efficient computation of $\text{PMSE}_{\text{CV}}$

Theorem (Theorem 2.7, Fan, Li, Zhang, and Zou 2020)

For any linear smoother  $\hat{\mathbf{Y}} = \mathbf{W}\mathbf{Y}$  with the self-stability property, we have

$$y_i - \hat{m}_{(i)}(x_i) = \frac{y_i - \hat{y}_i}{1 - w_{ii}}$$

and, therefore, its leave-one-out cross validation error is

$$PMSE_{CV} = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{y}_i}{1 - w_{ii}} \right)^2.$$

## Proof for the ridge regression estimation

- Let  $\hat{\beta}_{\text{ridge},\lambda}^{(i)}$  be the estimation of  $\beta = (\beta_1, \dots, \beta_p)$  when leaving out the  $i$ -th observation:

$$\hat{\beta}_{\text{ridge},\lambda}^{(i)} = \arg \min_{\beta} \left\{ \sum_{l=1, l \neq i}^n \left( y_l - \sum_{j=1}^p x_{lj} \beta_j \right)^2 + \lambda \|\beta\|_2^2 \right\}$$

- Let us define

$$\tilde{y}_l^{(i)} = \begin{cases} y_l & \text{if } l \neq i, \\ \hat{y}_i^{(i)} = \sum_{j=1}^p x_{ij} \hat{\beta}_{\text{ridge},\lambda,j}^{(i)} & \text{if } l = i. \end{cases}$$

- It follows that for all  $\beta \in \mathbb{R}^p$ ,

$$\sum_{l=1}^n \left( \tilde{y}_l^{(i)} - \sum_{j=1}^p x_{lj} \beta_j \right)^2 + \lambda \|\boldsymbol{\beta}\|_2^2 =$$

$$\left\{ \sum_{l=1, l \neq i}^n \left( y_l - \sum_{j=1}^p x_{lj} \beta_j \right)^2 + \lambda \|\boldsymbol{\beta}\|_2^2 \right\} + \left( \sum_{j=1}^p x_{ij} (\hat{\beta}_{\text{ridge}, \lambda, j}^{(i)} - \beta_j) \right)^2$$

- Observe that  $\hat{\beta}_{\text{ridge},\lambda}^{(i)}$  minimizes both terms in the right hand side. Then it is also

$$\hat{\beta}_{\text{ridge},\lambda}^{(i)} = \arg \min_{\beta} \left\{ \sum_{l=1}^n \left( \tilde{y}_l^{(i)} - \sum_{j=1}^p x_{lj} \beta_j \right)^2 + \lambda \|\beta\|_2^2 \right\}$$

- This is the ridge regression estimator corresponding to a data set with matrix of explanatory variables  $\mathbf{X}$  and vector of responses  $\tilde{\mathbf{y}}^{(i)} = (\tilde{y}_1^{(i)}, \dots, \tilde{y}_n^{(i)})^\top$ .
- Then

$$\begin{aligned} \hat{\beta}_{\text{ridge},\lambda}^{(i)} &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \tilde{\mathbf{y}}^{(i)}, \\ \hat{\mathbf{y}}^{(i)} &= \mathbf{X} \hat{\beta}_{\text{ridge},\lambda}^{(i)} = \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \tilde{\mathbf{y}}^{(i)} = \mathbf{H}_\lambda \tilde{\mathbf{y}}^{(i)}. \end{aligned}$$



- Observe that the  $i$ -th element of  $\hat{\mathbf{y}}^{(i)}$  is just  $\hat{y}_i^{(i)} = \sum_{j=1}^p x_{ij} \hat{\beta}_{\text{ridge}, \lambda, j}^{(i)}$ .
- Let  $\mathbf{e}_i$  be the  $n$ -dimensional vector whose  $i$ -th element is 1 and the others are equal to 0.
- Then  $\tilde{\mathbf{y}}^{(i)} = \mathbf{y} - (y_i - \hat{y}_i^{(i)}) \mathbf{e}_i$  and, consequently,

$$\hat{\tilde{\mathbf{y}}}^{(i)} = \mathbf{H}_\lambda \tilde{\mathbf{y}}^{(i)} = \mathbf{H}_\lambda \left( \mathbf{y} - (y_i - \hat{y}_i^{(i)}) \mathbf{e}_i \right) =$$

$$\mathbf{H}_\lambda \mathbf{y} - (y_i - \hat{y}_j^{(i)}) \mathbf{H}_\lambda \mathbf{e}_i = \hat{\mathbf{y}} - (y_i - \hat{y}_j^{(i)}) \mathbf{h}_j^\lambda,$$

where  $\mathbf{h}_i^\lambda$  is the  $i$ -th column of  $\mathbf{H}_\lambda$ .

- Looking just at the  $i$ -th component,  $\hat{y}_i^{(i)} = \hat{y}_i - (y_i - \hat{y}_i^{(i)})h_{ii}^\lambda$ , where  $h_{ii}^\lambda$  is the element  $(i, i)$  of  $\mathbf{H}_\lambda$ , or the  $i$ -th element in the diagonal of  $\mathbf{H}_\lambda$ .
- Then  $y_i - \hat{y}_i^{(i)} = y_i - \hat{y}_i + (y_i - \hat{y}_i^{(i)})h_{ii}^\lambda$  and we conclude that

$$y_i - \hat{y}_i^{(i)} = \frac{y_i - \hat{y}_i}{1 - h_{ii}^\lambda}.$$

So the loo-CV errors  $(y_i - \hat{y}_i^{(i)})$  can be computed if we know the errors  $(y_i - \hat{y}_i)$  when fitting the ridge regression with all the data, and the diagonal of  $\mathbf{H}_\lambda$ , and the proof concludes.

## Practice:

- Prostate data: Efficient computation of  $\text{PMSE}_{\text{CV}}$  in ridge regression.
- Use the R script `prostate.ridge.regression.R`.

## Generalized cross-validation

- For **linear estimators** of the regression function, a modification can be done in the measure of  $PMSE_{CV}$ .
- It is known as **generalized cross-validation (GCV)**.
- It consists in replacing in the expression of  $PMSE_{CV}$  the values  $w_{ii}$ , coming from the diagonal of  $\mathbf{W}$ , by their average value:

$$\text{PMSE}_{\text{GCV}} = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{y}_i}{1 - \nu/n} \right)^2,$$

$\nu = \text{Trace}(\mathbf{W}) = \sum_{i=1}^n w_{ii}$  is the effective number of parameters.

- In ridge regression,  $\lambda_{\text{GCV}} = \arg \min_{\lambda} \text{PMSE}_{\text{GCV}}(\lambda)$ .
- Manipulating the expression of  $\text{PMSE}_{\text{GCV}}$  it follows that

$$\text{PMSE}_{\text{GCV}} = \frac{n\hat{\sigma}_{\varepsilon}^2}{n - \nu},$$

where  $\hat{\sigma}_\varepsilon^2 = \frac{1}{n-\nu} \sum_{i=1}^n (y_i - \hat{y}_i)^2$  estimates the residual variance.

## Practice:

- Prostate data:  $\text{PMSE}_{\text{GCV}}$  in ridge regression.
- Use the R script  
`prostate.ridge.regression.R`.

Fan, J., R. Li, C.-H. Zhang, and H. Zou (2020).

*Statistical foundations of data science.*

Chapman and Hall/CRC.

Hastie, T., R. Tibshirani, and J. Friedman (2009).

*The Elements of Statistical Learning* (2nd ed.).

Springer.

Hastie, T., R. Tibshirani, and M. Wainwright (2015).

### Statistical learning with sparsity: the lasso and generalizations.

CRC Press.