

Linear models

(Some basic results)

N. Torelli, L. Egidi, G. Di Credico

Fall 2020

University of Trieste

Matrix notation

Inference in Linear models

Model validation and model selection

Matrix notation

The linear model

- Linear models (LM) are appropriate when analyzing the relationship between a quantitative *response variable* Y and a set of *covariates* x_1, x_2, \dots, x_{p-1} .

It is assumed that a sample of n values of the response variable Y is observed as well as n values of each covariate.

- The aim is to evaluate the impact of covariates on the mean μ_i of the response variable Y_i for the i -th unit. In a linear model this is represented by the equation

$$E(Y_i) = \mu_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{p-1} x_{ip-1} \quad (1)$$

The value y_i for the i -th unit of the sample can be written as

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{p-1} x_{ip-1} + \epsilon_i, \quad (2)$$

the model above can be also written for the set of all the n units in matrix notation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- $\mathbf{X}\boldsymbol{\beta}$ is the so called systematic component

Matrix notation

The equation for each unit in matrix notation is

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p-1} \\ \vdots & & & \vdots \\ 1 & x_{n1} & \cdots & x_{np-1} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_{p-1} \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix},$$

con:



- $\mapsto \mathbf{y}$ = is the vector of the values of the response variable ($n \times 1$) ;
- $\mapsto \mathbf{X}$ = is a matrix ($n \times p$) which contains the values of the covariates.
- $\mapsto \boldsymbol{\beta}$ = is the vector ($p \times 1$) of the regression coefficients;
- $\mapsto \boldsymbol{\epsilon}$ = is the vector ($n \times 1$) of the stochastic components.

The model written for the i -th unit can be also written as $y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i$
where \mathbf{x}_i^T is the i -th row of the *so-called* design matrix.

Matrix notation: main assumptions

In the linear model

- The response variable Y is a quantitative variable
- The covariates X could be either:
 - ↳ quantitative (numeric) variables or
 - ↳ categorical variables (factors).
- it is usually assumed that the values in the matrix \mathbf{X} are fixed constant (non stochastic). \mathbf{X} is the **design matrix**
- the design matrix \mathbf{X} is assumed to be of full rank. Since usually $n \gg p$ this means that the rank of \mathbf{X} is p (that is $\min(p, n)$).

The columns of \mathbf{X} are linearly independent vectors.

Matrix notation: main assumptions

The linear model is completely specified by assumptions on the stochastic components, the random variables ϵ_i

1. $E(\epsilon_i) = 0$ or equivalently $E(\epsilon) = 0$
2. $Var(\epsilon_i) = \sigma^2$ homoscedasticity
3. $E(\epsilon_i, \epsilon_j) = 0$ per $i \neq j$ uncorrelation.

The last two conditions can be more concisely expressed in matrix form as

$$Cov(\epsilon) = E(\epsilon\epsilon^T) = \sigma^2 I_n,$$

where $Cov(\epsilon)$ denotes variance-covariance matrix of the random vector ϵ .

The assumption 1-3 are called the second order assumptions (since they refer only to the first two moments of the variables).

A distributional assumption is then often added

4. $\epsilon \sim N_n(0, \Sigma)$

In matrix notation

$$E(y) = X\beta \text{ and } y \sim N(X\beta, \sigma^2 I)$$

Discussion of the assumptions

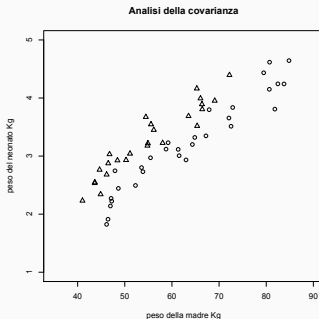
- **Linearity.** The assumption about the linear effect of the covariates is actually not very restrictive. **Non linear relationships can be introduced by appropriate transformations of the covariates.** For instance:
 - $y_i = \beta_0 + \beta_1 \log(z_i) + \epsilon_i$ introduces a logarithmic effect of z_i . But note that if one simply redefines $x_i = \log(z_i)$ then we are back to a standard linear model for the transformed variable X .
 - x_i^2 introduces a parabolic effect.
- **Homoscedasticity of the random components.** This is the standard assumption. The use of diagnostic checks can help verifying it. If possible departures from homoscedasticity are ignored it can impair quality of estimates. Possible remedies can be introduced
- **Uncorrelated random components.**
It is assumed that all random components are mutually uncorrelated. In some context this assumption is questionable (this is the case of data that are temporally or spatially ordered). Also this assumption can be verified with diagnostic tools and remedies are available.

Continuous covariates, factors, interactions

- When the covariate is a quantitative one, under the linearity assumption, then the value of the parameter associated to it represents simply the derivative of Y wrt to X .
- The effect of a categorical variable (a factor) measures the difference in the expected value of the response variable for each value of the factor wrt the reference category for the factor itself (all the other variables being equal). There will be then as many parameters as the number of levels of the factor minus one.
- Usually the interactions between two (or more) variables are introduced. Interpretation of interaction is easier when it refers to two factors or to a factor and a numeric variable.

Another example: one factor and a quantitative variable

Weight, Y , in kilograms, for a sample of newborn babies, from smoker mothers smokers (F) (in the graph different symbols are used for smokers - triangles - and non smokers - circles). For each women the pre-pregnancy weight mother weight X is observed)



As expected the weights of newborn babies is greater on the average when the mothers are non smokers. It seems that a systematic difference exists between the two groups though the relationship between weight of the mother and weight of the babies does not change.

A model with a continuous covariate and a factor

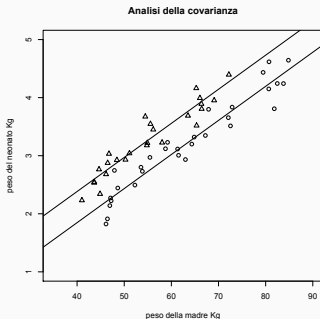
The two variables Y ="weight of the babies" and X ="weight of the mother" are both continuous numeric variables while the variable F is a factor with two levels ($F = 1$ if smoker, $F = 2$ if non smoker). The model is

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 I_{(F_i=2)} + \epsilon_i .$$

with the corresponding matrices:

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_1 & 0 \\ 1 & x_2 & 0 \\ \vdots & \vdots & \vdots \\ 1 & x_k & 1 \\ \vdots & \vdots & \vdots \\ 1 & x_n & 1 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix} .$$

Interpretation of the parameters



For the given specification

- parameter β_1 measures the (linear) effect of the weight of the mother on the weight of the baby and represents the (common) slope of the two lines in the graph;
- β_0 measures the intercept of the smoker's line and β_2 is, for a given weight of the mother, the vertical distance between the two lines.

An alternative parametrization

- Any model, particularly when factors are involved, can have alternative parametrizations.

The model introduced above can be also written in the following form:

$$Y_i = \beta_1 x_i + \beta_2 I_{(F_i=1)} + \beta_3 I_{(F_i=2)} + \epsilon_i, \quad (3)$$

where $I_{(F_i=j)}$ is an indicator variable which takes on 1 if $(F_i = j)$, $j = 1, 2$, and 0 otherwise.

- The model is equivalent but interpretation of parameters changes.

In this case the matrix form is:

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} x_1 & 1 & 0 \\ x_2 & 1 & 0 \\ \vdots & \vdots & \vdots \\ x_n & 0 & 1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}.$$

Interpretation of the parameters

In this new parametrization

- parameter β_1 measures the effect of the weight of the mother on the weight of the baby
- β_2 and β_3 estimate the mean of the y the weight of the babies, for a given weight of the mother, for smokers and non smokers respectively.
- The columns of the design matrix \mathbf{X} can be obtained from those of the previous specification by using a linear combination. The model is the same but interpretation of parameters changes.
- one could not add the intercept otherwise \mathbf{X} will become rank deficient

A model with interaction

The independent variables enter the model additively: the effect of the variable X is the same for each level of the factor F . In many case this is a too simplistic model, and the effect can change for different values of the factor.

This effect can be caught by introducing interaction between the covariates. The following regression model includes an interaction

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 I_{(F_i=2)} + \beta_3 I_{(F_i=2)} x_i + \epsilon_i .$$

Interaction implies the relationship between x (weight of the mother) and y (weight of the baby) can be different for smokers and non smokers.

Inference in Linear models

Least square estimation

- A sample y_1, y_2, \dots, y_n along with the values of the vector \mathbf{x}_i for the covariates observed on i -th unit is available and the aim is to estimate the parameters of the models $(\beta_0, \beta_1, \dots, \beta_{p-1}, \sigma^2)$
- An estimator of the parameters vector β can be obtained by using the **least square method**:

1. The **least square estimator (LSE)** $\hat{\beta}$ of β is the vector for which the following quantity is minimized

$$LS(\beta) = \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) ;$$

$$LS(\beta) = \mathbf{y}^T \mathbf{y} - \beta^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \beta + \beta^T \mathbf{X}^T \mathbf{X} \beta$$

2. Taking the derivative $\frac{\partial LS(\beta)}{\partial \beta} = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \beta$ and then equating it to 0, the LSE is

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

To invert $(\mathbf{X}^T \mathbf{X})$ we have to assume that this matrix is **non singular**.
This is **always true if \mathbf{X} is a full rank matrix**

Properties of LS estimator

The LSE $\hat{\beta}$ has the following properties:

1. $E(\hat{\beta}) = \beta$ and $\text{var}(\hat{\beta}) = \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$;
2. asymptotically $\hat{\beta} \sim N_p(\beta, \sigma^2 V^{-1})$ where $V = \lim_{n \rightarrow \infty} \mathbf{X}_n^T \mathbf{X}_n$ and \mathbf{X}_n is the sequence of design matrices and V is positive definite; in most of the cases then for large n , $\hat{\beta} \sim N_p(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$. This allows us to test significance of parameters and to build confidence intervals easily.
3. $\hat{\beta}$, the LSE is the *best estimator* in the sense that it has minimum variance among all linear estimators (BLUE Best Linear Unbiased Estimator - Gauss-Markov theorem).
4. When $(\mathbf{X}^T \mathbf{X})$ is not singular but its determinant is very close to 0 then estimates are very unstable. This happens if (multiple) correlation among the column of \mathbf{X} is very close to 1. Regularization could be a solution.

- When normality is assumed for the random components and taking account of uncorrelation (which means also independence in the normal case) then β can be obtained by using *maximum likelihood estimation*.
- Choose the value β which maximizes

$$L(\beta, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta)\right)$$

It is easy to show (by evaluating the log-likelihood and by deriving with respect to β) that the likelihood is maximized if the following quantity is minimized $(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta)$.

- **under normality assumption, MLE and LSE are equivalent**
- But to the properties already listed for the LSE now it can be added the one regarding the (exact) distribution of $\hat{\beta} \sim N(\beta, \sigma^2(\mathbf{X}^T\mathbf{X})^{-1})$

Estimation of σ^2

- To estimate σ^2 the following estimator is usually considered

$$\begin{aligned} S^2 &= \frac{SSE}{n - p} \quad \text{where} \\ SSE &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2 \\ &= (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) \\ &= (\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta}) \end{aligned}$$

is the *Sum of Squares of residuals* e_i .

- Note also that, when normality holds

$$\hat{\beta} \sim N_p(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}) \quad \text{and} \quad \frac{(n - p)S^2}{\sigma^2} \sim \chi_{n-p}^2,$$

and the two estimator $\hat{\beta}$ and S^2 are independent.

- Note that for small samples if σ^2 is unknown and S^2 is used, tests and confidence intervals for a single β_j are based on student t distribution with $n - p$ df.

Model validation and model selection

Testing a general linear hypothesis

The tests of more common interest are:

- test of significance for a single element of β

$$H_0 : \beta_j = 0 \quad \text{against} \quad H_1 : \beta_j \neq 0$$

- test on a subvector $\beta_1 = (\beta_1, \dots, \beta_r)$

$$H_0 : \beta_1 = 0 \quad \text{against} \quad H_1 : \beta_1 \neq 0$$

- test of equality of two coefficients $H_0 : \beta_j - \beta_r = 0$ against

$$H_1 : \beta_j - \beta_r \neq 0$$

All these hypotheses are special cases of the *general linear hypothesis*

$$H_0 : \mathbf{C}\beta = \mathbf{d} \quad \text{against} \quad H_1 : \mathbf{C}\beta \neq \mathbf{d}$$

\mathbf{C} is a $r \times p$ matrix with $\text{rank} = r \leq p$ and \mathbf{d} is a $r \times 1$ vector.

If data are fit to the model under the restriction $\mathbf{C}\beta = \mathbf{d}$ the residuals of this model are $_{H_0}e_i$ and one can compute $SSE_{H_0} = \sum_i^n _{H_0}e_i^2$ and calculate the statistic

$$\frac{n-p}{r} \frac{SSE_{H_0} - SSE}{SSE}$$

that, when H_0 is true, under normality assumption is a $F_{r,n-p}$ random variable.

Test significance of a single coefficient

When testing significance for a single element of β

$$H_0 : \beta_j = 0 \quad \text{against} \quad H_1 : \beta_j \neq 0$$

applying general linear hypothesis (where C is a row vector of zeros with one only in position j and $d = 0$), it can be shown that

$$\frac{\hat{\beta}_j^2}{\text{Var}(\hat{\beta}_j)} \sim F_{1, n-p}$$

This is the square of a student t with $n - p$ df and equivalently the following test statistic is usually considered

$$t_j = \frac{\hat{\beta}_j}{\text{Var}(\hat{\beta}_j)^{1/2}}$$

This result can be also used to obtain for β_j the following confidence interval at level $1 - \alpha$:

$$\beta_j \pm t_{n-p, 1-\alpha/2} (\text{Var}(\hat{\beta}_j))^{1/2}$$

Decomposition of Sum of Squares

- The following holds:

$$SST = SSR + SSE, \quad (4)$$

the **Total Sum of Squares** (*total deviance*) is the sum of the **Regression Sum of Squares** (*deviance explained by the model*) and the **Residual Sum of squares** (*deviance of the residuals*). Analyzing the components of (4) is of great relevance, **the ratio between SSR and SST is** clearly **related to the quality of the model.**

- Let \mathcal{F}_1 be the minimal model (the one which contains only the intercept, $p = 1$).

Let \mathcal{F}_p be the current model with p parameters and let

\mathcal{F}_{p_o} be a reduced model with $1 < p_o < p$ nested in \mathcal{F}_p .

Then the variance explained by the current model \mathcal{F}_p can be partitioned as it is shown in the table that follows (Table 1), called **Analysis of variance table.**

The analysis of variance

Table 1: Analysis of Variance (Anova)

Source of variability	df	SS	testing models improvement
total	n	SST	
constant	1	$n\bar{Y}^2$	
total	$n - 1$	SST_{cor}	
improvement with \mathcal{F}_{p_o} with respect to \mathcal{F}_1	$p_o - 1$	SSR_{p_o}	$\frac{SSR_{p_o} / (p_o - 1)}{SSE_{p_o} / (n - p_o)}$ $\sim F_{p_o - 1, n - p_o}$
improvement with \mathcal{F}_p with respect to \mathcal{F}_{p_o}	$p - p_o$	$SSR_p - SSR_{p_o}$	$\frac{(SSE_{p_o} - SSE_p) / (p - p_o)}{SSE_p / (n - p)}$ $\sim F_{p - p_o, n - p}$
residuals \mathcal{F}_p	$n - p$	SSE_p	

- The fall in the fit from \mathcal{F}_{p_o} to \mathcal{F}_p can be evaluated using the statistic

$$F = \frac{(SSE_{p_o} - SSE_p) / (p - p_o)}{SSE_p / (n - p)} \sim F_{p - p_o, n - p}.$$

Coefficient of determination R^2

The coefficient of determination R^2 is defined as the proportion of total variance explained by the regression model.

It can be used as a goodness-of-fit measure for the models

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{SSE}{SST}$$

and $0 \leq R^2 \leq 1$.

This decomposition is possible if the model includes the intercept.

For nested models R^2 always increases adding covariates.

When comparing nested models the corrected coefficient of determination \bar{R}^2 is instead used

$$\bar{R}^2 = 1 - \frac{n-1}{n-p}(1 - R^2)$$

It penalizes inclusion of new variables that are non significant.

- The mean of \mathbf{y} can be predicted once the model is estimated by $\hat{\mu}_i = \hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ and consequently $\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = H\mathbf{y}$
- $H = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ is a square matrix of size n and it is called the **hat matrix** or the projection matrix. It has the following properties:
 1. it is symmetric and idempotent
 2. $\text{rank}(H) = \text{trace}(H) = p$
 3. h_{ii} have values that range from $1/n$ and 1 and their sum is equal to p
 4. the matrix $I - H$ is also symmetrical and idempotent with rank equal to $(n - p)$
- The residuals of the model are then $\mathbf{e}_i = (I - H)\mathbf{y}$
- under normality assumption $\mathbf{e}_i \sim N(\mathbf{0}, \sigma^2(I - H))$

Analysis of the residuals

- The quality of the model and the validity of the assumptions can be judged by using some diagnostic tools that mainly rely upon analysis of residuals as defined above.
- In the linear models the residuals can be **standardized** (to take into account the fact that they have unequal variance)

$$r_i = \frac{Y_i - \hat{\mu}_i}{\sqrt{S^2(1 - h_{ii})}} , \quad (5)$$

where h_{ii} is the i -th element on the diagonal of $H = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$. These values are called **leverages**. They reveal if a point has values that are far from the majority of data in the space of the x s. Suspect values have leverage $> 2p/n$

- to **identify which values are outliers** with respect the majority of the data points the **studentized residual** are introduced

$$e_i^* = \frac{e_i}{S_{(i)}\sqrt{1 - h_{ii}}}$$

where $S_{(i)}^2$ is the variance of the residuals when the i -th observation is excluded. For a model correctly specified e_i^* follows a t with $n - p - 1$ df.

Analysis of the residuals

- Cook distances are defined as

$$D_i = \frac{1}{p} r_i^2 \frac{h_{ii}}{1 - h_{ii}} = \frac{1}{pS^2} \sum_{j=1}^n (\hat{y}_j - \hat{y}_{j(i)})^2$$

- $\hat{y}_{j(i)}$ are the predicted values for the i -th observation in a model estimated without it. It reveals (when $D_i > 1$) observations that, when excluded from the analysis, will cause substantial modifications in the estimates of the parameters.
- Classical graphical tools based on (standardized) residuals are:
 - plot of residuals against the predicted values (*to reveal possible heteroscedasticity*)
 - plot of residuals against the explanatory variables (*to reveal non linearities*)
 - plot of residuals against variables not in the model (*added variable plot*)
 - Q-Q norm of residuals (*to assess normality*)
 - plot of leverages h_{ii} and of Cook distances (*to reveal outliers*)
- formal test of normality, such as Shapiro-Wilks test or Jarque-Bera test (the latter based on estimated values of third and fourth standardized moments)

Dealing with non constant variance and residual correlation

- Residual analysis can reveal that some assumptions could be questionable
- Critical assumptions are those on uncorrelation and heteroscedasticity of residuals
- The assumptions $Cov(\epsilon) = \sigma^2 I$ should be replaced by a more general $Cov(\epsilon) = \sigma^2 W^{-1}$ where W is assumed to be simply a positive definite matrix
- In this case LSE is still an unbiased estimator but the estimate of its variance covariance matrix is biased.
- If we ignore this, the main consequences are that tests or confidence intervals based on assumption of uncorrelation and homoscedasticity lead to wrong conclusions (tests tends to say a parameter is significant too often and confidence intervals appear shorter)

Heteroscedasticity

- We will consider here only the case of heteroscedasticity. In this case $\text{Cov}(\epsilon) = \sigma^2 W^{-1}$ with $W^{-1} = \text{diag}(1/w_1, 1/w_2, \dots, 1/w_n)$
- If each ϵ_i is multiplied by $\sqrt{w_i}$ one obtains the transformed values $\epsilon_i^* = \sqrt{w_i}\epsilon_i$ which have constant variance
- $\text{var}(\epsilon_i^*) = \text{var}(\sqrt{w_i}\epsilon_i) = \sigma^2$ and then the random components are homoscedastic.
- The model does not change if we transform also the response variable and the covariates (including the intercept) accordingly.
- We then obtain

$$y_i^* = \sqrt{w_i}y_i \text{ and}$$

$$x_{ij}^* = \sqrt{w_i}x_{ij}$$

for each of the p covariates (including the intercept) and then the model

$$y_i^* = \beta_0 x_{i0}^* + \beta_1 x_{i1}^* + \beta_2 x_{i2}^* + \dots + \beta_{p-1} x_{i,p-1}^* + \epsilon_i^*$$

is homoscedastic and the same assumptions of a standard LM hold. These transformations, in matrix notation, are equivalent to pre-multiply all components of the model $\mathbf{y} = \mathbf{X}\beta + \epsilon$ by the matrix $W^{1/2}$

Weighed Least Squares (WLS)

- If then

$$\mathbf{y}^* = W^{1/2}\mathbf{y}, \mathbf{X}^* = W^{1/2}\mathbf{X} \text{ and } \epsilon^* = W^{1/2}\epsilon$$

we get a new model for transformed data where homoscedasticity holds and parameters can again be estimated by LSE obtaining

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}^{*T}\mathbf{X}^*)^{-1}\mathbf{X}^{*T}\mathbf{y}^* \\ &= (\mathbf{X}^T W^{1/2} W^{1/2} \mathbf{X})^{-1} \mathbf{X}^T W^{1/2} W^{1/2} \mathbf{y} \\ &= (\mathbf{X}^T W \mathbf{X})^{-1} \mathbf{X}^T W \mathbf{y}\end{aligned}$$

- This estimator is the **Weighted Least Square Estimator** (WLS)
- Note that weights are inversely proportional to variances of ϵ (which are originally heteroscedastic)
- To units with a more erratic random component are given smaller weights
- **Application of this strategy requires that the weights are known**

Model choice and variable selection

In many applications a large number of candidate predictors are available.
A naive approach often used is the following:

Estimate the most complex model that includes all the covariates (and possibly all the interactions). Then, remove all not significant variables from the model. (Backward selection)

This strategy is not advisable for many reasons. Let us list some of them:

- the resulting model can overfit the data and then its predictive performance for new data can decrease
- the larger the number of covariate the higher the risk of multi-collinearity (correlated regressors)
- There are many models with equivalent performances but different substantial interpretation. You are not sure that the variable which remains in your model after such backward selection strategy should be really considered the most relevant.

Other naive (yet often used) strategies for variable selection are:

- All subset selection (chose the best among $\sum_{j=1}^P \binom{P}{j}$ possible models)
- Forward selection
- Stepwise selection (a combination of forward and backward selection)

Model choice criteria

Since one of the principles to consider when building a model is the *Occam's razor*, criteria to select a model that has good performances and at the same time is less complex should be introduced.

when considering alternative LM we have already seen some criteria

- R^2 and corrected R^2
- F test (for nested models)
- Mallows's C_p

$$C_p = \frac{\sum_i^n (y_i - \hat{y}_{iM})^2}{\hat{\sigma}^2} - n + M$$

where M is the number of covariates in the model and \hat{y}_{iM} are the predicted values with those M covariates. The "best" model is the one with lowest C_p .

- Akaike Information Criteria (AIC)

$$AIC = -2l(\hat{\beta}_M, \hat{\sigma}^2) + 2(M + 1)$$

Better fit corresponds to smaller smaller AIC values.

For a linear model with gaussian components and p β_j parameters

$$AIC = n \log(\hat{\sigma}^2) + 2(p + 1).$$

Note that $\hat{\sigma}^2$ is SSE divided by n .

- Bayesian Information Criteria (BIC)

$$BIC = -2l(\hat{\beta}_M, \hat{\sigma}^2) + \log(n)(M + 1)$$

Avoiding collinearity

A diagnosis of collinearity is obtained by computing the variance inflation factor (VIF) associated to the j -th predictor

$$VIF_j = \frac{1}{1 - R_j^2}$$

where R_j^2 is the coefficient of determination when x_j is regressed on all the remaining covariates. $VIF_j > 10$ is usually taken as a symptom that the variable can cause collinearity.

Typical solutions are:

- omission of covariates
- using principal components extracted from regressors (or other combination of the regressors)
- Ridge regression: it is an alternative to LSE where

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

and λ is a chosen tuning parameter

Regularization Techniques

It has been assumed so far that \mathbf{X} has full rank, and this gives a unique solution to equations

$$(\mathbf{X}^T \mathbf{X})\beta = \mathbf{X}^T \mathbf{y}$$

We have already discussed that $\mathbf{X}^T \mathbf{X}$ could be close to singularity. Regularization consists in changing the objective function by penalizing it:

$$\hat{\beta}_{PLS} = \arg \min_{\beta} [(\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \text{pen}(\beta)]$$

where $\text{pen}(\beta)$ is a term that measure the complexity of the model and $\lambda \geq 0$ is a smoothing parameter that reflects the weight given to the penalty. Penalty is such that it is large when many β are large.

Ridge regression is an example of penalized least square. It corresponds to introducing the following penalty

$$\text{pen}(\beta) = \sum_{j=1}^p \beta_j^2 = \beta^T \beta$$

With large λ the penalty term dominates and all (or almost all) coefficients are shrunk to 0. λ is usually chosen by k -fold cross validation.

LASSO (Least Absolute Shrinkage and Selection Operator)

Also LASSO corresponds to a penalized least square criterion and

$$\hat{\beta}_{LASSO} = \arg \min_{\beta} [(\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \sum_{j=1}^p |\beta_j|]$$

- The penalization chosen with LASSO tend to shrink some of the values of the coefficients to 0. Small coefficients are more strongly shrunk to 0 compared with Ridge regression.
- Balance between fit of the data and regularization
- Note that no closed explicit solution of the minimization problem exists. Numerical optimization must be used (quadratic programming, LARS–Least Angle Regression).