SF 2930 REGRESSION ANALYSIS LECTURE 10

Diagnostics & Validation: PRESS, Cook's Distance

Timo Koski

KTH Royal Institute of Technology

2023

LEARNING OUTCOMES

- Crossvalidation
- Crossvalidation by PRESS =predicted residual error sum of squares
 - LOOCV = leave-one-out-crossvalidation
 - Sherman-Morrison-Woodbury (S-M-W) theorem and formula: the problem of determining the change in the inverse matrix after a change in a column.
- Expectation of PRESS
- PRESS in Model Choice, Hjorth's Theorem.
- Cook's Distance for Detecting Data Points that are Influential for LSE $\hat{\beta}$.
- Recursive Multiple Regression

CROSSVALIDATION

Cross-validation or out-of-sample testing, is any of several similar model validation techniques for assessing how the results of a statistical analysis will generalize to an independent data set. Cross-validation is a resampling method that uses different portions of the data to test and train a model on different iterations.

It is mainly used in settings where the goal is prediction, and one wants to estimate how accurately a predictive model will perform in practice.

CROSSVALIDATION

The calculations of a summary measure of the fit of a model on the test set, a sample of observations that were not themselves used to estimate the model. This is a simple idea and easy to implement. But it has a few potential drawbacks

- Test error rate can be highly variable, depending on precisely which observations are included in the training set and which observations are included in the test set.
- Statistical methods tend to perform worse when trained on fewer observations (which is a result splitting data into tow sets).

Methods of crossvalidation (CV) are developed to overcome this problem, and do model validation by resampling without splitting the data.

LOOCV= LEAVE ONE OUT CROSSVALIDATION

The predicted residual error sum of squares (PRESS) is one of techniques of cross-validation used in linear regression analysis. A fitted model having been produced, each observation in turn is removed and the model is refitted using the remaining observations. The out-of-sample predicted value is calculated for the omitted observation in each case, and the PRESS statistic is calculated as the sum of the squares of all the resulting prediction errors.

CROSSVALIDATION BY PRESS

The training set is

$$\mathcal{D}_{tr} = \left\{ \left(\mathbf{x}_{i}^{\mathsf{T}}, y_{i}\right)_{j=1}^{n} \right\}.$$

We remove $(\mathbf{x}_i^{\mathsf{T}}, y_i)$ and the new training set is

$$\mathcal{D}_{tr,(i)} = \left\{ \left(\mathbf{x}_{i}^{\mathsf{T}}, y_{i} \right)_{j=1, j \neq i}^{n} \right\}$$

One finds the LSE $\widehat{\boldsymbol{\beta}}_{(i)}$ using $\mathcal{D}_{tr,(i)}$. Then the predictor of y_i is $\widehat{\boldsymbol{y}}_{(i)} = \mathbf{x}_{(i)}^{\mathsf{T}} \widehat{\boldsymbol{\beta}}_{(i)}$. We repeat this is for all $(\mathbf{x}_i^{\mathsf{T}}, y_i)$, $i = 1, \ldots, n$. The PRESS statistic is

PRESS =
$$\sum_{i=1}^{n} (y_i - \hat{y}_{(i)})^2$$



CROSSVALIDATION BY PRESS. MODEL CHOICE

The PRESS statistic is

PRESS =
$$\sum_{i=1}^{n} (y_i - \hat{y}_{(i)})^2$$

Given this procedure, the PRESS statistic can be calculated for a number of candidate model structures for the same dataset, with the lowest values of PRESS indicating the best structures.

Models that are over-parameterised (over-fitted) would tend to give small residuals for observations included in the model-fitting but large residuals for observations that are excluded.

CROSSVALIDATION BY PRESS

This is straighforward as an idea, but is this practical, as it seems that we must fit n multiple regressions to data.



However, there is a short way: only one LSE based on \mathcal{D}_{tr} is needed.

REFRESHER

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

$$\mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_i \\ \vdots \\ \mathbf{y}_n \end{pmatrix}, X = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_i^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{i1} & \cdots & x_{ik} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{pmatrix}$$

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}.$$

LOOCV: REMOVE $(\mathbf{x}_i^{\mathsf{T}}, y_i)$

$$\mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_l \\ \vdots \\ \mathbf{y}_n \end{pmatrix}, X = \begin{pmatrix} \mathbf{x}_1^\mathsf{T} \\ \mathbf{x}_2^\mathsf{T} \\ \vdots \\ \mathbf{x}_n^\mathsf{T} \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{l1} & \cdots & x_{lk} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{pmatrix}$$

THE (ORDINARY) MULTIPLE LINEAR REGRESSION MODEL

$$\beta \in \mathbb{R}^{k+1}$$
 and $n > k+1$.

$$\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}. \tag{1}$$

The following assumptions hold:

- 1) $E[\varepsilon] = \mathbf{0} \in \mathbb{R}^n$
- 2) $C_{\varepsilon} = E\left[\varepsilon \varepsilon^{T}\right] = \sigma^{2} \mathbb{I}_{n}$ (homoscedasticity)
- 3) X has full column rank



LSE

$$\mathcal{D}_{tr} = \{ \left(\mathbf{x}_{i}^{\mathsf{T}}, y_{i} \right)_{i=1}^{n} \}$$

We use the training set to compute

$$\widehat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y} = \operatorname{argmin}_{\boldsymbol{\beta}} \parallel \mathbf{y} - X \boldsymbol{\beta} \parallel^2$$

HAT MATRIX

$$H = X(X^T X)^{-1} X^T. (2)$$

The predicted values are $\hat{\mathbf{y}} = X \hat{\boldsymbol{\beta}} = X (X^T X)^{-1} X^T \mathbf{y} = H \mathbf{y}$. The elements on the main diagonal of H, h_{ii} , are called the leverages of the respective data point.

The $n \times 1$ random vector $\hat{\varepsilon}$ of observed residuals

$$\widehat{\varepsilon} = \mathbf{Y} - X\widehat{\boldsymbol{\beta}} = \mathbf{Y} - H\mathbf{Y}. \tag{3}$$

We have found that the covariance matrix of $\hat{\varepsilon}$ is

$$C_{\widehat{\varepsilon}} = \sigma^2 \left(\mathbb{I}_n - H \right) \tag{4}$$



THE SHORT WAY

PRESS =
$$\sum_{i=1}^{n} (y_i - \hat{y}_{(i)})^2$$

We shall establish that

$$\hat{y}_i = h_{ii} y_i + (1 - h_{ii}) \hat{y}_{(i)}, \tag{5}$$

where h_{ii} is the element on the main diagonal of H.



THE SHORT WAY

From (5)

$$y_i - \hat{y}_i = (1 - h_{ii})y_i - (1 - h_{ii})\hat{y}_{(i)},$$

which gives

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

Hence

$$PRESS = \sum_{i=1}^{n} \left(\frac{y_i - \hat{y}_i}{1 - h_{ii}} \right)^2.$$
 (6)

This requires, of course, that $h_{ii} < 1$, which is checked in an Appendix. This is a consequence of idempotency of H. MPV on p. 134, p. 151, p. 592 does not pay attention to this detail.

EXAM GENERATOR

Mahalanobis distance d_M^2 between \mathbf{x}_i and the arithmetic mean $\bar{\mathbf{x}}$ of the rows of X (the first 1 excluded). S is the empirical covariance matrix of \mathbf{x}_i .

$$d_M^2(\mathbf{x}_i, \bar{\mathbf{x}}; S) = (\mathbf{x}_i - \bar{\mathbf{x}})^T S^{-1} (\mathbf{x} - \bar{\mathbf{x}}).$$

The (possible) exam question is to show that

$$(n-1)\left(h_{ii}-\frac{1}{n}\right)=O_M^2(\mathbf{x}_i,\bar{\mathbf{x}};S).$$

Thus

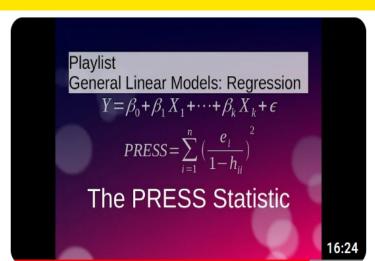
$$h_{ii} = \frac{1}{n-1} d_M^2(\mathbf{x}_i, \bar{\mathbf{x}}; S) + \frac{1}{n}$$

and then the PRESS statistic becomes

PRESS =
$$\sum_{i=1}^{n} \left(\frac{y_i - \hat{y}_i}{1 - \frac{1}{n-1} d_M^2(\mathbf{x}_i, \bar{\mathbf{x}}; S) - \frac{1}{n}} \right)^2$$
 (7)

YOUTUBE

https://www.youtube.com/watch?v=SBYO0dPbA



We prove

$$\hat{y}_i = h_{ii} y_i + (1 - h_{ii}) \hat{y}_{(i)}, \tag{8}$$

by deriving the following.

$$\widehat{\boldsymbol{\beta}}_{(i)} = \widehat{\boldsymbol{\beta}} + \frac{(\widehat{y}_i - y_i)}{1 - h_{ii}} (X^T X)^{-1} \mathbf{x}_i.$$
 (9)



SIMPLE STEPS ON THE SHORT WAY AFTER REMOVING $(\mathbf{x}_{i}^{\mathsf{T}}, y_{i})$: STEP ONE

$$X^{T}X = \begin{pmatrix} \mathbf{x}_{1}^{T} \\ \mathbf{x}_{2}^{T} \\ \vdots \\ \mathbf{x}_{n}^{T} \end{pmatrix} \begin{pmatrix} \mathbf{x}_{1}^{T} \\ \mathbf{x}_{2}^{T} \\ \vdots \\ \mathbf{x}_{n}^{T} \end{pmatrix} = (\mathbf{x}_{1}\mathbf{x}_{2} \dots \mathbf{x}_{n}) \begin{pmatrix} \mathbf{x}_{1}^{T} \\ \mathbf{x}_{2}^{T} \\ \vdots \\ \mathbf{x}_{n}^{T} \end{pmatrix} = \sum_{j=1}^{n} \mathbf{x}_{j} \mathbf{x}_{j}^{T}$$

$$= \sum_{j=1, j \neq i}^{n} \mathbf{x}_{j} \mathbf{x}_{j}^{T} + \mathbf{x}_{i} \mathbf{x}_{i}^{T} = X_{(i)}^{T} X_{(i)} + \mathbf{x}_{i} \mathbf{x}_{i}^{T}$$
(10)

$$X^{T}\mathbf{y} = (\mathbf{x}_{1}\mathbf{x}_{2}...\mathbf{x}_{n})\begin{pmatrix} y_{1} \\ y_{2} \\ \vdots \\ y_{n} \end{pmatrix} = X_{(i)}^{T}\mathbf{y}_{(i)} + \mathbf{x}_{i}y_{i}$$
(11)

: STEP TWO

From (11)

$$\widehat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y} = (X^T X)^{-1} X_{(i)}^T \mathbf{y}_{(i)} + (X^T X)^{-1} \mathbf{x}_i y_i$$
 (12)

From (10)

$$\widehat{\boldsymbol{\beta}}_{(i)} = (X_{(i)}^{T} X_{(i)})^{-1} X_{(i)}^{T} \mathbf{y} = \left(X^{T} X - \mathbf{x}_{i} \mathbf{x}_{i}^{T} \right)^{-1} X_{(i)}^{T} y_{i}.$$
 (13)

This is the *i*th regression due to removal of (\mathbf{x}_i^T, y_i) , but will not be computed along the short way.



STEP 3

Sherman-Morrison-Woodbury (S-M-W) theorem and formula: the problem of determining the change in the inverse matrix after a change in a column. In the current situation we are determining the change in the inverse matrix X^TX , when one row has been removed in X.

S-M-W

Suppose A is is an invertible square $n \times n$ matrix and $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ are column vectors. Then $A + \mathbf{u}\mathbf{v}^T$ is invertible iff $1 + \mathbf{v}^T A^{-1}\mathbf{u} \neq 0$. In this case,

$$(A + \mathbf{u}\mathbf{v}^{\mathsf{T}})^{-1} = A^{-1} - \frac{A^{-1}\mathbf{u}\mathbf{v}^{\mathsf{T}}A^{-1}}{1 + \mathbf{v}^{\mathsf{T}}A^{-1}\mathbf{u}}.$$
 (14)

We apply this with $A = X^T X$, $\mathbf{u} \mathbf{v}^T = -\mathbf{x}_i \mathbf{x}_i^T$.



SHERMAN-MORRISON-WOODBURY THEOREM

We apply S-M-W Formula (14) with $A = X^TX$, $\mathbf{u}\mathbf{v}^T = -\mathbf{x}_i\mathbf{x}_i^T$. We note that X^TX is invertible, since $X_{(i)}$ has full column rank, a standing assumption under n > k+1. Then we must check that $1 + \mathbf{v}^TA^{-1}\mathbf{u} \neq 0$, i.e.,

$$1-\mathbf{x}_i^{\mathsf{T}}(X^{\mathsf{T}}X)^{-1}\mathbf{x}_i\neq 0.$$

This detail is omitted/overlooked in MVP p. 591. However, $h_{ii} = \mathbf{x}_i^\mathsf{T} X^\mathsf{T} X \mathbf{x}_i$ is the *i*th element on the main diagonal of the hat matrix

$$H = X(X^T X)^{-1} X^T$$

It is shown in an Appendix that $h_{ii} < 1$.



SHERMAN-MORRISON-WOODBURY THEOREM

$$(X^{T}X - \mathbf{x}_{i}\mathbf{x}_{i}^{T})^{-1}X_{(i)}^{T}y_{i} = (X^{T}X)^{-1}X_{(i)}^{T}y_{i}$$
$$-\frac{(X^{T}X)^{-1}\mathbf{x}_{i}\mathbf{x}_{i}^{T}(X^{T}X)^{-1}}{1 - h_{ii}}X_{(i)}^{T}y_{i}.$$

Thus

$$\widehat{\boldsymbol{\beta}}_{(i)} = \underbrace{\left(\boldsymbol{X}^{T}\boldsymbol{X}\right)^{-1}\boldsymbol{X}_{(i)}^{T}\boldsymbol{y}_{i}}_{\equiv A} - \underbrace{\frac{\left(\boldsymbol{X}^{T}\boldsymbol{X}\right)^{-1}\boldsymbol{x}_{i}\boldsymbol{x}_{i}^{T}\left(\boldsymbol{X}^{T}\boldsymbol{X}\right)^{-1}}{1-h_{ii}}\boldsymbol{X}_{(i)}^{T}\boldsymbol{y}_{i}}_{\equiv B}.$$
 (15)

Now we re-write the right hand side of (15) by means of (12), that is, by insertion of

$$A = (X^{T}X)^{-1}X_{(i)}^{T}\mathbf{y}_{(i)} = \widehat{\beta} - (X^{T}X)^{-1}\mathbf{x}_{i}y_{i}.$$
 (16)

We substitute A in B to get

$$B = \frac{(X^{T}X)^{-1} \mathbf{x}_{i} \mathbf{x}_{i}^{T} (X^{T}X)^{-1}}{1 - h_{ii}} X_{(i)}^{T} y_{i} = \frac{(X^{T}X)^{-1} \mathbf{x}_{i} \mathbf{x}_{i}^{T} (\widehat{\beta} - (X^{T}X)^{-1} \mathbf{x}_{i} y_{i})}{1 - h_{ii}}$$

Here

$$(X^{T}X)^{-1} \mathbf{x}_{i} \mathbf{x}_{i}^{T} (\widehat{\boldsymbol{\beta}} - (X^{T}X)^{-1} \mathbf{x}_{i} y_{i}) = (X^{T}X)^{-1} \mathbf{x}_{i} \underbrace{\mathbf{x}_{i}^{T} \widehat{\boldsymbol{\beta}}}_{=\widehat{y}_{i}}$$
$$- (X^{T}X)^{-1} \mathbf{x}_{i} \underbrace{\mathbf{x}_{i}^{T} (X^{T}X)^{-1} \mathbf{x}_{i}}_{y_{i}} y_{i}$$

Hence

$$B = \frac{\left(X^{T}X\right)^{-1}\mathbf{x}_{i}\left(\widehat{y}_{i} - h_{ii}y_{i}\right)}{1 - h_{ii}}.$$
(17)

FINAL STEPS

When we insert (16) and (17) in (15) in we get

$$\widehat{\boldsymbol{\beta}}_{(i)} = A - B
= \widehat{\boldsymbol{\beta}} - (X^T X)^{-1} \mathbf{x}_i y_i
- \frac{(X^T X)^{-1} \mathbf{x}_i (\widehat{y}_i - h_{ii} y_i)}{1 - h_{ii}}$$

and this equals

$$= \widehat{\beta} + \frac{(X^{T}X)^{-1} \mathbf{x}_{i} \widehat{y}_{i}}{1 - h_{ii}} + \frac{-(X^{T}X)^{-1} \mathbf{x}_{i} (1 - h_{ii}) y_{i} - (X^{T}X)^{-1} \mathbf{x}_{i} h_{ii} y_{i}}{1 - h_{ii}}$$

$$= \widehat{\beta} + \frac{(X^{T}X)^{-1} \mathbf{x}_{i} \widehat{y}_{i}}{1 - h_{ii}} - \frac{(X^{T}X)^{-1} \mathbf{x}_{i} y_{i}}{1 - h_{ii}}$$

$$= \widehat{\beta} + \frac{(\widehat{y}_{i} - y_{i})}{1 - h_{ii}} (X^{T}X)^{-1} \mathbf{x}_{i}.$$

We have thus shown how to compute $\widehat{\beta}_{(i)}$ by of one LSE.

$$\widehat{\boldsymbol{\beta}}_{(i)} = \widehat{\boldsymbol{\beta}} + \frac{(\widehat{\mathbf{y}}_i - \mathbf{y}_i)}{1 - h_{ii}} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i.$$
 (18)

Next we establish (8).

$$\hat{\mathbf{y}}_{(i)} = \mathbf{x}_i^T \widehat{\boldsymbol{\beta}}_{(i)}.$$

We use (18) to find

$$\hat{\mathbf{y}}_{(i)} = \mathbf{x}_i^T \widehat{\boldsymbol{\beta}} + \frac{(\widehat{\mathbf{y}}_i - \mathbf{y}_i)}{1 - h_{ii}} \mathbf{x}_i^T (X^T X)^{-1} \mathbf{x}_i.$$

Here $\mathbf{x}_{i}^{T}\widehat{\boldsymbol{\beta}} = \widehat{y}_{i}$ and $h_{ii} = \mathbf{x}_{i}^{T}(X^{T}X)^{-1}\mathbf{x}_{i}$. Hence

$$\hat{\mathbf{y}}_{(i)} = \hat{\mathbf{y}}_i + \frac{\mathbf{h}_{ii} \left(\hat{\mathbf{y}}_i - \mathbf{y}_i \right)}{1 - \mathbf{h}_{ii}}$$

and

$$(1-h_{ii})\widehat{y}_{(i)}=(1-h_{ii})\widehat{y}_i+h_{ii}\left(\widehat{y}_i-y_i\right).$$

After little simplification this becomes (8).



PROPERTIES OF PRESS

In view of (3) we have

PRESS =
$$\sum_{i=1}^{n} \left(\frac{y_i - \hat{y}_i}{1 - h_{ii}} \right)^2 = \sum_{i=1}^{n} \frac{\hat{\epsilon}_i^2}{(1 - h_{ii})^2}$$

Therefore

$$E[PRESS] = \sum_{i=1}^{n} \frac{E\left[\hat{\epsilon}_{i}^{2}\right]}{(1 - h_{ii})^{2}}$$

By (4) we have that $E\left[\hat{\epsilon}_{i}^{2}\right]=\sigma^{2}(1-h_{ii})$. Hence

$$E[PRESS] = \sigma^2 \sum_{i=1}^{n} \frac{1}{(1 - h_{ii})^2}$$

$$PRESS = \sum_{i=1}^{n} \frac{\hat{\epsilon}_{i}^{2}}{(1 - h_{ii})^{2}}$$

Interpretation as follows: for a fair estimate of the error in predicting y_i , we should inflate the square training error $\hat{\epsilon}_i^2$ by a factor of $\frac{1}{(1-h_{ii})^2}$. Hence, the closer the value of h_{ii} to 1, the more we should inflate the training error. This makes sense when we recall (see Exam questions) that h_{ii} is the leverage of the ith data point, which is a measure of how much our prediction \hat{y}_i depends on the actual value y_i . The more $hat y_i$ depends on y_i , the more \hat{y}_i is an overfit, and the greater we should inflate the error squared error $\hat{\epsilon}_i^2$.

Data points with large residuals (outliers) and/or high leverage may distort the outcome and accuracy of a regression. Cook's distance, denoted by D_i , measures the effect of deleting a given observation. Hence this involves an application of (18). Points with a large Cook's distance are considered to merit closer examination in the analysis. Cook (1977) introduced D_i by

$$D_{i} := \frac{\left(\widehat{\boldsymbol{\beta}}_{(i)} - \widehat{\boldsymbol{\beta}}\right)^{T} \boldsymbol{X}^{T} \boldsymbol{X} \left(\widehat{\boldsymbol{\beta}}_{(i)} - \widehat{\boldsymbol{\beta}}\right)}{(k+1)\widehat{\sigma}^{2}}$$
(19)

Here we have the familiar $\hat{\sigma}^2 = \frac{\mathbf{e}^{\mathsf{T}}\mathbf{e}}{n-k-1}$. Next Cook substitutes for $\hat{\boldsymbol{\beta}}_{(i)} - \hat{\boldsymbol{\beta}}$ by (18) and obtains

$$D_i = \left(\frac{(\widehat{y}_i - y_i)}{\sqrt{(1 - h_{ii})}\widehat{\sigma}}\right)^2 \frac{\mathbf{x}_i^T (X^T X)^{-1} \mathbf{x}_i}{(k+1)(1 - h_{ii})}$$

But, see slides 25-29 of Lecture 4 or MPV p. 99, $\widehat{y}_i = \mathbf{x}_i^{\mathsf{T}} \widehat{\boldsymbol{\beta}}$ is the minimal MSE predictor of y_i and that $\mathrm{Var}\left[\widehat{y}_i\right] = \sigma^2 \mathbf{x}_i^{\mathsf{T}} \left(X^{\mathsf{T}} X\right)^{-1} \mathbf{x}_i$. In addition, by (4), $\mathrm{Var}\left[\widehat{\varepsilon}_i\right] = \sigma^2 = (1 - h_{ii}) = (1 - \mathbf{x}_i^{\mathsf{T}} \left(X^{\mathsf{T}} X\right)^{-1} \mathbf{x}_i)$. Thus we write

$$D_{i} = \left(\frac{(\widehat{y}_{i} - y_{i})}{\sqrt{(1 - h_{ii})}\widehat{\sigma}}\right)^{2} \frac{\mathbf{x}_{i}^{T}(X^{T}X)^{-1}\mathbf{x}_{i}}{(k+1)(1 - h_{ii})} = \frac{t_{i}^{2}}{(k+1)} \frac{\operatorname{Var}\left[\widehat{y}_{i}\right]}{\operatorname{Var}\left[\widehat{\varepsilon}_{i}\right]}$$

where

$$t_i = \frac{(\widehat{y}_i - y_i)}{\sqrt{(1 - h_{ii})}}, \quad i = 2, \dots, n$$

Note that $Var[t_i] = 1$ for every i.



$$D_i = \frac{t_i^2}{(k+1)} {(20)}$$

- t_i² is a measure of the degree to which the *i*th observation can be considered as an outlier from the assumed model.
 Any point *i* with with a large ŷ_i y_i and a large h_{ii} is potentially highly influential on the least squares fit.
- The ratios $\operatorname{Var}\left[\widehat{y}_{i}\right]/\operatorname{Var}\left[\widehat{\varepsilon}_{i}\right]$ measure the relative sensitivity of the estimate, $\widehat{\boldsymbol{\beta}}$, to potential outlying values at each data point.
- The use of D_i is possible only if $\operatorname{Var}\left[\widehat{\varepsilon}_i\right] > 0$.



Cooke (1977) discusses a regression analysis relating six economic variables to total derived employment for the years 1947 to 1962 in USA. Cook points out first that there are considerable differences between residuals e_i and t_i . Second, the point with the largest D_i value corresponds to 1951. In the table $R_i/s \leftrightarrow \widehat{\varepsilon}_i/\widehat{\sigma}$, t_i , $V(Y_i)/V(R_i)$ $\leftrightarrow Var[\widehat{y_i}]/Var[\widehat{\varepsilon}_i]$.

TABLE 1-Longley Data

Year	R _i /s	t _i	$V(\hat{Y}_i)/V(R_i)$	Di
1947	0.88	1.15	0.74	0.14
48	-0.31	0.48	1,30	0.04
49	0.15	0.19	0.57	*
50	-1.34	1.70	0.59	0.24
	1.02	1.64	1.60	0.61
51 52	-0.82	1.03	0.59	0.09
53	-0.54	0.75	0.97	0.08
54	-0.04	0.06	1.02	*
55	0.05	0.07	0.84	*
56	1.48	1.83	0.49	0.23
57	-0.06	0.07	0.56	*
58	-0.13	0.18	0.93	*
59	-0.51	0.64	0.60	0.04
60	-0.28	0.32	0.30	*
61	1.12	1.42	0.59	0.17
62	-0.68	1.21	2.21	0.47

*: smaller than 5 x 10^{-3}

TECHNOMETRICS©, VOL. 19, NO. 1, FEBRUARY 1977

STUDENTIZED RESIDUALS

The quantities in Cook's Distance

$$t_i = \frac{(\widehat{y}_i - y_i)}{\sqrt{(1 - h_{ii})}} = \frac{(\widehat{\epsilon}_i)}{\sqrt{(1 - h_{ii})}}, \quad i = 2, \dots, n$$

are a form of a Student's t-statistic, are therefore known as the studentized residuals, see MVP p. 133.

The key reason for studentizing in regression analysis is that the variances of the residuals at different data points may differ, even if the variances of the errors at these data points are equal.

$$\operatorname{Var}\left[\epsilon_{i}\right] = \sigma^{2}, \quad \operatorname{Var}\left[\widehat{\epsilon}_{i}\right] = (1 - h_{ii}),$$

where h_{ii} is the leverage of *i*th data point.



PRESS AND MODEL CHOICE

Models \mathcal{M}_k , $k = 1, 2, \dots, M$.

PRESS
$$(\mathcal{M}_k) = \sum_{i=1}^{n} \left(\frac{y_i - \hat{y}_i(\mathcal{M}_k)}{1 - h_{ii}(\mathcal{M}_k)} \right)^2$$

Choose \mathcal{M}_{k_*} , if

$$k_* = \operatorname{argmin}_{k \in \{1, \dots, M\}} \mathsf{PRESS}\left(\mathcal{M}_k\right)$$

URBAN HJORTH'S BIAS THEOREM

PROPOSITION

Models \mathcal{M}_k , $k = 1, 2, \dots, M$.

$$\mathcal{M}_{k_*} = \min_{k \in \{1, \dots, M\}} \mathsf{PRESS}\left(\mathcal{M}_k\right)$$

Then

$$E\left[\mathsf{PRESS}\left(\mathcal{M}_{k_*}\right)\right] \le \min_{k \in \{1, \dots, M\}} E\left[\mathsf{PRESS}\left(\mathcal{M}_k\right)\right]. \tag{21}$$

Proof: Define the r.v. $X_k := \text{PRESS}\left(\mathcal{M}_k\right) - \text{PRESS}\left(\mathcal{M}_{k_*}\right)$ for every outcome of X, \mathbf{Y} . Then $X_k \geq 0$. A result in probability calculus says that

If
$$X \ge 0$$
, and $P(X = 0) < 1$, then $E[X] > 0$.

Hence $E[X_k] > 0$, and (21) follows.



URBAN HJORTH'S BIAS THEOREM

The proposition above is on p. 35 of J.S. Urban Hjorth: *Computer Intensive Statistical Methods*, Chapman & Hall, London, 1994.

It says that there is a bias that depends on the model selection method. The expected predictive value of any single PRESS (\mathcal{M}_k) is larger, than the expected predictive value of the best model. Hence the best chosen model underestimates the expected predictive value of any PRESS (\mathcal{M}_k) , i.e., the best models give a too optimistic view.

RECURSIVE ESTIMATION FOR MULTIPLE LINEAR Models

A SECOND APPLICATION OF S-M-W: RECURSIVE ESTIMATION ALGORITHMS FOR LSE

$$\mathcal{D}_{tr,(i)} = \left\{ \left(\mathbf{x}_{i}^{\mathsf{T}}, y_{i}\right)_{j=1, j \neq i}^{n} \right\}$$

New data $(\mathbf{x}_{n+1}, y_{n+1})$ is received and the training set is augmented by it. The least-squares estimate $\widehat{\boldsymbol{\beta}}$ for the parameters is updated to reflect the new data. By S-M-V formula we show that the updated LSE $\widehat{\boldsymbol{\beta}}_{n+1}$ is given recursively as

$$\widehat{\boldsymbol{\beta}}_{n+1} = \widehat{\boldsymbol{\beta}} + C\left((\mathbf{x}_{n+1}, y_{n+1}), (X^T X)^{-1}\right)$$

where $C\left((\mathbf{x}_{n+1},y_{n+1}),(X^TX)^{-1}\right)$ is a $(k+1)\times 1$) the random vector that depends only on $(\mathbf{x}_{n+1},y_{n+1})$ and the current estimate of σ^2 and a storing of current $(X^TX)^{-1}$. Hence, we do not need to solve the normal equations again to get the updated LSE.

New data \mathbf{x}_{n+1} , y_{n+1} is received and the training set is augmented by it.

$$\mathbf{y}_{n+1} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_n \\ y_{n+1} \end{pmatrix}, X_{n+1} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_i^T \\ \vdots \\ \mathbf{x}_{n+1}^T \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{i1} & \cdots & x_{ik} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \\ 1 & x_{(n+1)1} & \cdots & x_{(n+1)k} \end{pmatrix}$$

The current desing matrix *X* is found inside the augmented matrix as

$$X_{n+1} = \begin{pmatrix} X \\ \mathbf{x}_{n+1}^T \end{pmatrix}$$
.

Hence we get

$$X_{n+1}^{T}X_{n+1} = \begin{pmatrix} X \\ \mathbf{x}_{n+1}^{T} \end{pmatrix}^{T} \begin{pmatrix} X \\ \mathbf{x}_{n+1}^{T} \end{pmatrix}$$

$$= \begin{pmatrix} X^{T} & \mathbf{x}_{n+1} \end{pmatrix} \begin{pmatrix} X \\ \mathbf{x}_{n+1}^{T} \end{pmatrix} = X^{T}X + \mathbf{x}_{n+1}\mathbf{x}_{n+1}^{T}.$$
(22)

As new data $(\mathbf{x}_{n+1}, y_{n+1})$ is received, the least-squares estimate for the parameters is updated to reflect the new data. The up-dated LSE $\widehat{\boldsymbol{\beta}}_{n+1}$ is

$$\widehat{\boldsymbol{\beta}}_{n+1} = \left(X_{n+1}^{\mathsf{T}} X_{n+1} \right)^{-1} X_{n+1}^{\mathsf{T}} \mathbf{y}_{n+1}$$
 (23)

Now can apply the S-H-W formula (14) and (22) for the inversion

$$(X_{n+1}^{T}X_{n+1})^{-1} = (X^{T}X + \mathbf{x}_{n+1}\mathbf{x}_{n+1}^{T})^{-1}$$

$$= (X^{T}X)^{-1} - \frac{(X^{T}X)^{-1}\mathbf{x}_{n+1}\mathbf{x}_{n+1}^{T}(X^{T}X)^{-1}}{1 + \mathbf{x}_{n+1}^{T}(X^{T}X)^{-1}\mathbf{x}_{n+1}}.$$
(24)

Let us note that $1 + \mathbf{x}_{n+1}^{\mathsf{T}} \left(X^{\mathsf{T}} X \right)^{-1} \mathbf{x}_{n+1} \neq 0$, since $\mathbf{x}_{n+1}^{\mathsf{T}} \left(X^{\mathsf{T}} X \right)^{-1} \mathbf{x}_{n+1} \geq 0$, because $\left(X^{\mathsf{T}} X \right)^{-1}$ is positive definite.

In other words, if the old inverse of (X^TX) is available, then the new inverse of $X_{n+1}^TX_{n+1}$ is found by adding a rank 1 correction to (X^TX) . In addition

$$X_{n+1}^{T} \mathbf{y}_{n+1} = \begin{pmatrix} X \\ \mathbf{x}_{n+1}^{T} \end{pmatrix}^{T} \begin{pmatrix} \mathbf{y} \\ y_{n+1} \end{pmatrix}$$

$$= \begin{pmatrix} X^{T} & \mathbf{x}_{n+1} \end{pmatrix} \begin{pmatrix} \mathbf{y} \\ y_{n+1} \end{pmatrix} = X^{T} \mathbf{y} + \mathbf{x}_{n+1} y_{n+1}.$$
(25)

46 / 52

We start now by using the S-H-W formula (24) in (23)

$$\widehat{\boldsymbol{\beta}}_{n+1} = \left(X_{n+1}^{T} X_{n+1} \right)^{-1} X_{n+1}^{T} \mathbf{y}_{n+1} =$$

$$= \left(X^{T} X \right)^{-1} X_{n+1}^{T} \mathbf{y}_{n+1} - \frac{\left(X^{T} X \right)^{-1} \mathbf{x}_{n+1} \mathbf{x}_{n+1}^{T} \left(X^{T} X \right)^{-1}}{1 + \mathbf{x}_{n+1}^{T} \left(X^{T} X \right)^{-1} \mathbf{x}_{n+1}} X_{n+1}^{T} \mathbf{y}_{n+1}$$

In order to simplify writing (actually $K = K_n$ and $\alpha = \alpha_n$), let us set

$$K := (X^T X)^{-1} \mathbf{x}_{n+1}, \quad \alpha := \frac{1}{1 + \mathbf{x}_{n+1}^T (X^T X)^{-1} \mathbf{x}_{n+1}}$$

Thus

$$\widehat{\boldsymbol{\beta}}_{n+1} = \left(\boldsymbol{X}^{T}\boldsymbol{X}\right)^{-1} \boldsymbol{X}_{n+1}^{T} \boldsymbol{y}_{n+1} - \alpha \boldsymbol{K} \boldsymbol{x}_{n+1}^{T} \left(\boldsymbol{X}^{T}\boldsymbol{X}\right)^{-1} \boldsymbol{X}_{n+1}^{T} \boldsymbol{y}_{n+1}$$

Now (25) in the first term in the right hand side entails

$$=\underbrace{\left(X^{T}X\right)^{-1}X^{T}\mathbf{y}}_{\widehat{\boldsymbol{x}}}+\underbrace{\left(X^{T}X\right)^{-1}\mathbf{x}_{n+1}}_{-K}y_{n+1}-\alpha K\mathbf{x}_{n+1}^{T}\left(X^{T}X\right)^{-1}X_{n+1}^{T}\mathbf{y}_{n+1}$$

Now we have come to

$$\widehat{\boldsymbol{\beta}}_{n+1} = \widehat{\boldsymbol{\beta}} + y_{n+1} \boldsymbol{K} - \alpha \boldsymbol{K} \boldsymbol{X}_{n+1}^{\mathsf{T}} \left(\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X} \right)^{-1} \boldsymbol{X}_{n+1}^{\mathsf{T}} \boldsymbol{y}_{n+1}. \tag{26}$$

We use again (25) to get

$$\alpha K \mathbf{x}_{n+1}^{\mathsf{T}} \left(X^{\mathsf{T}} X \right)^{-1} X_{n+1}^{\mathsf{T}} \mathbf{y}_{n+1} = \alpha K \mathbf{x}_{n+1}^{\mathsf{T}} \underbrace{\left(X^{\mathsf{T}} X \right)^{-1} X^{\mathsf{T}} \mathbf{y}}_{\widehat{\boldsymbol{\beta}}}$$

$$+\alpha K\mathbf{x}_{n+1}^{\mathsf{T}}\left(X^{\mathsf{T}}X\right)^{-1}\mathbf{x}_{n+1}y_{n+1}$$

Again by slides 25-29 of Lecture 4 or MPV p. 99 we get that $\widehat{y}_{n+1} = \mathbf{x}_{n+1}^{\mathsf{T}} \widehat{\boldsymbol{\beta}}$ is the minimal MSE predictor of y_{n+1} and that $\operatorname{Var}\left[\widehat{y}_{n+1}\right] = \sigma^2 \mathbf{x}_{n+1}^{\mathsf{T}} \left(X^{\mathsf{T}} X\right)^{-1} \mathbf{x}_{n+1}$.



When we insert in (26) we find

$$\widehat{\boldsymbol{\beta}}_{n+1} = \widehat{\boldsymbol{\beta}} + y_{n+1} K$$

$$-\alpha \widehat{y}_{n+1} K - \frac{\alpha}{\sigma^2} \operatorname{Var} \left[\widehat{y}_{n+1} \right] K.$$
(27)

Here

$$\frac{\alpha}{\sigma^2} = \frac{1}{\sigma^2 \left(1 + \mathbf{x}_{n+1}^{\mathsf{T}} \left(X^{\mathsf{T}} X\right)^{-1} \mathbf{x}_{n+1}\right)} = \frac{1}{\sigma^2 \left(1 + \frac{\operatorname{Var}\left[\widehat{y}_{n+1}\right]}{\sigma^2}\right)}$$
$$= \frac{1}{\sigma^2 + \operatorname{Var}\left[\widehat{y}_{n+1}\right]}$$

When we insert in (27) we get

$$\widehat{\boldsymbol{\beta}}_{n+1} = \widehat{\boldsymbol{\beta}} + c(n+1)K_{n+1}, \tag{28}$$

where the scalar c(n+1) is

$$c(n+1) = y_{n+1} - \frac{\sigma^2}{\sigma^2 + \operatorname{Var}\left[\widehat{y}_{n+1}\right]} \widehat{y}_{n+1} - \frac{1}{\sigma^2 + \operatorname{Var}\left[\widehat{y}_{n+1}\right]} \operatorname{Var}\left[\widehat{y}_{n+1}\right]$$

and the $(k + 1) \times 1$ vector K_n is

$$K_{n+1} = \left(X^T X\right)^{-1} \mathbf{x}_{n+1}.$$

The idea for this recursion is from

Hager, William W: Updating the inverse of a matrix. *SIAM Review*, 31, 2, pp. 221–239, 1989

But the final recursion in (28) is different from the one found loc.cit..

SIAM =Society for Industrial and Applied Mathematics

APPENDIX: $h_{ii} < 1$

 h_{ii} is on the main diagonal of the hat matrix H. We take a look at the ith element on the main diagonal of H^2 . Let \mathbf{h}_{ii} denote column i in H. Then the ith element on the main diagonal of H^2 is found as

$$\mathbf{h}_i^{\mathsf{T}}\mathbf{h}_i = \sum_{j=1}^n h_{ij}^2.$$

But $H^2 = H$. Hence

$$h_{ii} = \sum_{j=1}^{n} h_{ij}^2 = \sum_{j=1, j \neq i}^{n} h_{ij}^2 + h_{ii}^2$$

. Since $\sum_{j=1, j\neq i}^{n} h_j^2 > 0$ we get

$$h_{ii} > h_{ii}^2 \Leftrightarrow 1 > h_{ii}$$
.

Here we could divide both sides by h_{ii} , since $h_{ii} > 1/n > 0$, see one of the problems in SF2930Examgenerator 2023.