

SF 2930 REGRESSION ANALYSIS

LECTURE 4

Multiple Linear Regression, Part 2

Timo Koski

KTH Royal Institute of Technology

2023

LEARNING OUTCOMES

- Gauss Markov Theorem
- Prediction
- Statistical Properties of the Residuals
- Distributions of Quadratic Forms
 - χ^2 -distribution

PART I: REFRESHMENT

Selected Findings from Preceding Lectures Required in this Lecture.

EXQ PYTHAGORAS'S THEOREM

Show that

$$\| \mathbf{y} \|^2 = \| \hat{\mathbf{y}} \|^2 + \| \mathbf{e}_{LSE} \|^2$$

The data point \mathbf{y} is the hypotenuse of the right-angled triangle in \mathbb{R}^n with the base of predicted/fitted values $\hat{\mathbf{y}}$ and the altitude of the LSE- residual \mathbf{e}_{LSE} . This is next illustrated in a Figure.

By Courtesy of Puntanen, S. and Isotalo, J. and Styan, GPH:
Formulas Useful for Linear Regression Analysis and Related Matrix Theory. In the Figure $\hat{\epsilon} \leftrightarrow \mathbf{e}_{LSE}$

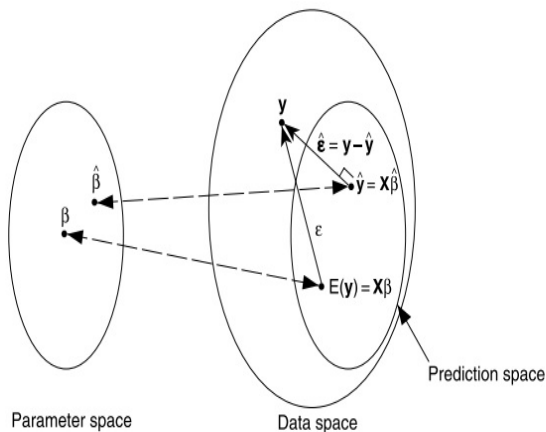


Figure 7.4 Geometric relationships of vectors associated with the multiple linear regression model.

PROJECTION GEOMETRICALLY FOR SIMPLE LINEAR REGRESSION

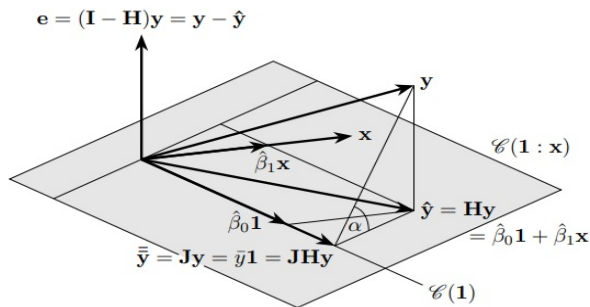


Figure 8.3 Projecting y onto $\mathcal{C}(1 : x)$.

From Puntanen S., Styan G.P.H., Isotalo J.: *Matrix Tricks for Linear Statistical Models*. Springer 2011.

LINEAR TRANSFORMATIONS

PROPOSITION

\mathbf{Y} and \mathbf{X} are random vectors, $\mu_{\mathbf{Y}} = E[\mathbf{Y}]$, $\mu_{\mathbf{X}} = E[\mathbf{X}]$, \mathbf{X} has covariance matrix $C_{\mathbf{X}}$, A and B are $m \times n$ matrices. \mathbf{a} and \mathbf{b} are vectors of suitable dimension. Then we have

- $E[\mathbf{X} + \mathbf{Y}] = \mu_{\mathbf{X}} + \mu_{\mathbf{Y}}$

- $\mathbf{Z} = A\mathbf{X} + \mathbf{b}$,

$$E[\mathbf{Z}] = A\mu_{\mathbf{X}} + \mathbf{b}, \quad (1)$$

$$C_{\mathbf{Z}} = AC_{\mathbf{X}}A^T. \quad (2)$$

- $C_{\mathbf{X}} = E[\mathbf{X}\mathbf{X}^T] - \mu_{\mathbf{X}}\mu_{\mathbf{X}}^T$

-

$$\text{Var}[\mathbf{a}^T \mathbf{X}] = \mathbf{a}^T C_{\mathbf{X}} \mathbf{a} \quad (3)$$

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

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, X = \begin{pmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{pmatrix}$$

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

THE (ORDINARY) MULTIPLE LINEAR REGRESSION MODEL. $k > 1$ COVARIATES/PREDICTORS

$$\beta \in \mathbb{R}^{k+1}, n \geq k+1.$$

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon. \quad (4)$$

The following assumptions hold:

- 1) $E[\epsilon] = \mathbf{0} \in \mathbb{R}^n$
- 2) $C_\epsilon = E[\epsilon\epsilon^T] = \sigma^2\mathbb{I}_n$ (homoscedasticity)
- 3) $X^T X$ is invertible

The model is called ordinary normal regression model, if additionally the following assumption holds:

- 4) $\epsilon \in N_n(\mathbf{0}, \sigma^2\mathbb{I}_n)$

LEAST SQUARES ESTIMATION

$$Q(\beta) := \| \mathbf{y} - X\beta \|^2 \quad (5)$$

and the LSE is the minimizer

$$\hat{\beta} := \operatorname{argmin}_{\beta} Q(\beta).$$

PROPOSITION

If $X^T X$ is a positive definite matrix, then

$$\hat{\beta} = (X^T X)^{-1} X^T \mathbf{y}. \quad (6)$$

HAT MATRIX

$$H := X(X^T X)^{-1} X^T. \quad (7)$$

$$\hat{\mathbf{y}} = H\mathbf{y} \in \text{sp}(X).$$

SUMMARY: ORDINARY MULTIPLE REGRESSION

$$\mathbf{Y} = X\boldsymbol{\beta}_* + \boldsymbol{\varepsilon}. \quad \text{True model} \quad (8)$$

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{Y}$$

$$E[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}_*$$

$$C_{\hat{\boldsymbol{\beta}}} = \sigma^2 (X^T X)^{-1} \quad (9)$$

$$\mathbf{e}_{LSE} = \mathbf{y} - X\hat{\boldsymbol{\beta}} = \mathbf{y} - H\mathbf{y}$$

$$\hat{\sigma}^2 = \frac{1}{(n - k - 1)} \mathbf{e}_{LSE}^T \mathbf{e}_{LSE}$$

SUMMARY: NORMAL (GAUSSIAN) MULTIPLE REGRESSION

$\varepsilon \in N(\mathbf{0}, \sigma^2 \mathbb{I}_n)$ and β_* such that

$$\mathbf{Y} = X\beta_* + \varepsilon \quad \text{True model} \quad (10)$$

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T \mathbf{Y} \\ \hat{\beta} &\sim N_{k+1}(\beta_*, \sigma^2 (X^T X)^{-1}) \end{aligned} \quad (11)$$

$$\begin{aligned} \mathbf{e}_{LSE} &= \mathbf{y} - X\hat{\beta} = \mathbf{y} - H\mathbf{y} \\ \hat{\sigma}^2 &= \frac{1}{(n - k - 1)} \mathbf{e}_{LSE}^T \mathbf{e}_{LSE} \end{aligned} \quad (12)$$

$$\widehat{\sigma^2} = \frac{1}{(n - k - 1)} \mathbf{e}_{LSE}^T \mathbf{e}_{LSE} \quad (13)$$

The χ^2 distribution of the quadratic form in the LSE residuals determining $\widehat{\sigma^2}$, $\mathbf{e}_{LSE}^T \mathbf{e}_{LSE}$, will be eventually derived for the normal multiple regression in this Lecture.

PART II

Gauss-Markov Theorem

Gauss-Markov theorem states that $\hat{\beta} = (X^T X)^{-1} X^T \mathbf{Y}$ has uniquely the lowest variance within the class of linear unbiased estimators

GAUSS-MARKOV THEOREM

PROPOSITION

Let $\tilde{\beta}$ be any unbiased linear estimator of β_* . Let $\hat{\beta} = (X^T X)^{-1} X^T \mathbf{Y}$, i.e. the unbiased ordinary LSE. Then it holds for an arbitrary $(k+1) \times 1$ vector \mathbf{a} that

$$\text{Var} \left[\mathbf{a}^T \tilde{\beta} \right] \geq \text{Var} \left[\mathbf{a}^T \hat{\beta} \right]. \quad (14)$$

If $\text{Var} \left[\mathbf{a}^T \tilde{\beta} \right] = \text{Var} \left[\mathbf{a}^T \hat{\beta} \right]$, then $\tilde{\beta} = \hat{\beta}$.

Proof: Let $\tilde{\beta} = B\mathbf{Y} + g_o$ be another unbiased linear estimator, where g_o is a $(k+1) \times 1$ vector. Unbiasedness means that $E[\tilde{\beta}] = \beta_*$. On the other hand, by the rule (1) and the true model (8) above

$$E[\tilde{\beta}] = E[B\mathbf{Y} + g_o] = BE[\mathbf{Y}] + g_o = BX\beta_* + g_o.$$

For unbiasedness it must hold that

$$BX = \mathbb{I}_{k+1}, g_o = \mathbf{0}_{k+1} \quad (15)$$

Now we take without loss of generality that

$$B = (X^T X)^{-1} X^T + G. \quad (16)$$

Then

$$BX = (X^T X)^{-1} X^T X + GX = \mathbb{I}_{k+1} + GX. \quad (17)$$

In view of (15) it holds that

$$GX = \mathbf{0}_{k+1}. \quad (18)$$

Next we find the covariance matrix of $\tilde{\beta}$. By the rule (2) and the true model (8)

$$C_{\tilde{\beta}} = BC_Y B^T = \sigma^2 BB^T.$$

Here

$$\begin{aligned} BB^T &= \left((X^T X)^{-1} X^T + G \right) \left((X^T X)^{-1} X^T + G \right)^T \\ &= \left((X^T X)^{-1} X^T + G \right) \left(X (X^T X)^{-1} + G^T \right) \\ &= (X^T X)^{-1} \underbrace{X^T X (X^T X)^{-1}}_{=I_{k+1}} + (X^T X)^{-1} \underbrace{X^T G^T}_{=\mathbf{0}_{k+1}^T} \\ &\quad + \underbrace{GX}_{=\mathbf{0}_{k+1}} (X^T X)^{-1} + GG^T, \end{aligned}$$

where we used (18) twice, since $X^T G^T = (GX)^T = \mathbf{0}_{k+1}^T$.

Thus

$$C_{\tilde{\beta}} = \sigma^2 \left((X^T X)^{-1} + G G^T \right) \quad (19)$$

Let now \mathbf{a} be an arbitrary $(k+1) \times 1$ vector. By (3) and (19)

$$\text{Var} \left[\mathbf{a}^T \tilde{\beta} \right] = \mathbf{a}^T C_{\tilde{\beta}} \mathbf{a} = \sigma^2 \mathbf{a}^T (X^T X)^{-1} \mathbf{a} + \sigma^2 \mathbf{a}^T G G^T \mathbf{a} \quad (20)$$

Put $\mathbf{z} = G^T \mathbf{a}$. Then $\mathbf{a}^T G G^T \mathbf{a} = \mathbf{z}^T \mathbf{z} = \|\mathbf{z}\|^2 \geq 0$. Hence

$$\text{Var} \left[\mathbf{a}^T \tilde{\beta} \right] \geq \sigma^2 \mathbf{a}^T (X^T X)^{-1} \mathbf{a}. \quad (21)$$

Due to (9) we have $\sigma^2 \mathbf{a}^T (X^T X)^{-1} \mathbf{a} = \mathbf{a}^T C_{\hat{\beta}} \mathbf{a}$. Again, by the rule (2) we have $\mathbf{a}^T C_{\hat{\beta}} \mathbf{a} = \text{Var} \left[\mathbf{a}^T \hat{\beta} \right]$. In other words, in (21)

$$\text{Var} \left[\mathbf{a}^T \tilde{\beta} \right] \geq \text{Var} \left[\mathbf{a}^T \hat{\beta} \right],$$

which is (14).

Next we prove uniqueness of $\hat{\beta}$ as stated in the theorem. We have found in (20)

$$\text{Var} [\mathbf{a}^T \tilde{\beta}] = \text{Var} [\mathbf{a}^T \hat{\beta}] + \sigma^2 \mathbf{a}^T G G^T \mathbf{a}$$

Hence if $\text{Var} [\mathbf{a}^T \tilde{\beta}] = \text{Var} [\mathbf{a}^T \hat{\beta}]$, we have for any \mathbf{a} the equality

$$\mathbf{a}^T G G^T \mathbf{a} = 0.$$

As above we set $\mathbf{z} = G^T \mathbf{a}$, and then $\mathbf{a}^T G G^T \mathbf{a} = \mathbf{z}^T \mathbf{z} = \|\mathbf{z}\|^2 = 0$. But a norm is point-separating, i.e., $\|\mathbf{z}\|^2 = 0$ implies that $\mathbf{z} = \mathbf{0}_{k+1}$.

Thus $G^T \mathbf{a} = \mathbf{0}_{k+1}$ for every \mathbf{a} . Hence, we are allowed to take \mathbf{a} as the standard basis vector $\mathbf{e}_j = (0, \dots, 0, \underbrace{1}_{\text{position } j}, 0, \dots, 1)^T$, so that

$G^T \mathbf{e}_j$ is the j :th row of G , which is $= \mathbf{0}_{k+1}^T$. In this way we recognize that every row in G is the zero vector $\mathbf{0}_{k+1}^T$.

Hence $G = \mathbf{0}$, = the $(k + 1) \times n$ zero matrix, and by (16) it follows that

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

Hence we have shown that $\text{Var} [\mathbf{a}^T \tilde{\boldsymbol{\beta}}] = \text{Var} [\mathbf{a}^T \hat{\boldsymbol{\beta}}]$, implies $\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}$, i.e., the asserted uniqueness property. \square

COROLLARY

For $j = 0, \dots, k$

$$\text{Var} [\tilde{\beta}_j] \geq \text{Var} [\hat{\beta}_j].$$

Proof: Use the standard basis vectors

$\mathbf{e}_j = (0, \dots, 0, \underbrace{1}_{\text{position } j}, 0, \dots, 1)^T$ in (14) for $j = 0, \dots, k$. \square

REMARK

The preceding proofs of the proposition and its corollary do not require the multivariate normal distribution, and the Gauss-Markov theorem is thus valid for any ordinary multiple LSE.

PREDICTION IN REAL TIME

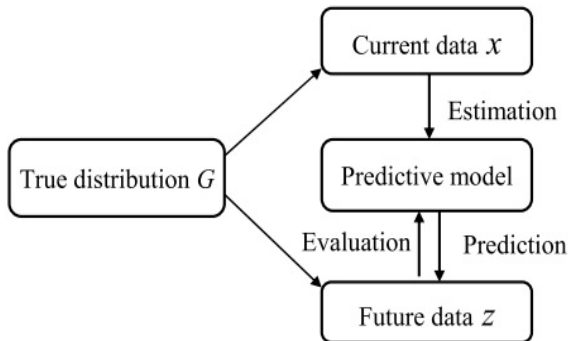


Fig. 1.2. Statistical modeling and the predictive point of view.

PREDICTION IN REAL TIME

A common situation is that we want to forecast a new value y_{n+1} based on the values of the \mathbf{x} -covariates. If we have the LSE $\hat{\beta}$ of β_* , then an unbiased (to be shown) prediction is

$$\hat{y}_{n+1} = \mathbf{x}_{n+1}^T \hat{\beta}, \quad \text{where } \mathbf{x}_{n+1}^T = (1, x_{n+1,1}, \dots, x_{n+1,k})$$

We can think of prediction in real time. We have observed the responses y_1, \dots, y_n , up to time n and wish to predict the next value, y_{n+1} . We assume, of course, that the underlying “true” mechanism generating data is unchanged in the sense that y_{n+1} is an outcome of

$$Y_{n+1} = \mathbf{x}_{n+1}^T \beta_* + \epsilon_{n+1}.$$

Note that ϵ_{n+1} is assumed to be independent of $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T$.

PREDICTION IN REAL TIME

In order to simplify writing (and to accomodate to other possible cases of prediction) we set

$$\hat{y} = \hat{y}_{n+1}, \mathbf{x} = \mathbf{x}_{n+1}, \quad e = \epsilon_{n+1}, Y_{n+1} = \mathbf{x}^T \boldsymbol{\beta}_* + e.$$

PREDICTION IN REAL TIME: ERROR

Let us proceed to calculate the prediction error $Y_{n+1} - \hat{y}$. We have

$$Y_{n+1} - \hat{y} = \mathbf{x}^T \boldsymbol{\beta}_* + e - \mathbf{x}^T \hat{\boldsymbol{\beta}} = \mathbf{x}^T (\boldsymbol{\beta}_* - \hat{\boldsymbol{\beta}}) + e$$

We see that the expected prediction error equals zero

$$E[Y_{n+1} - \hat{y}] = \mathbf{x}^T E[(\boldsymbol{\beta}_* - \hat{\boldsymbol{\beta}})] + E[e] = 0,$$

as $\hat{\boldsymbol{\beta}}$ is unbiased and $e \sim N(0, \sigma^2)$. Next apply

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T (X \boldsymbol{\beta}_* + \boldsymbol{\varepsilon}) = \boldsymbol{\beta}_* + (X^T X)^{-1} X^T \boldsymbol{\varepsilon}$$

to get

$$Y_{n+1} - \hat{y} = -\mathbf{x}^T (X^T X)^{-1} X^T \boldsymbol{\varepsilon} + e.$$

Note that $e \sim N(0, \sigma^2)$ and $\boldsymbol{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbb{I}_n)$.

MEAN SQUARE ERROR (MSE) OF PREDICTION IN REAL TIME

We strive to compute $MSE := E[(Y_{n+1} - \hat{y})^2]$. Here $\mathbf{x}^T(X^T X)^{-1}X^T \epsilon$ is a univariate r.v.. We square to get

$$(Y_{n+1} - \hat{y})^2 = \left(\mathbf{x}^T(X^T X)^{-1}X^T \epsilon\right)^2 - 2\left(\mathbf{x}^T(X^T X)^{-1}X^T \epsilon\right) \cdot e + e^2. \quad (22)$$

ϵ and e are independent, hence

$$E\left[\left(\mathbf{x}^T(X^T X)^{-1}X^T \epsilon\right) \cdot e\right] = \mathbf{x}^T(X^T X)^{-1}X^T E[\epsilon] \cdot E[e] = 0. \quad (23)$$

MEAN SQUARE ERROR (MSE) OF PREDICTION IN REAL TIME

Since $\mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}$ holds for scalar products

$$\begin{aligned} E \left[\left(\mathbf{x}^T (X^T X)^{-1} X^T \boldsymbol{\varepsilon} \right)^2 \right] &= E \left[\mathbf{x}^T (X^T X)^{-1} X^T \boldsymbol{\varepsilon} \cdot \mathbf{x}^T (X^T X)^{-1} X^T \boldsymbol{\varepsilon} \right] \\ &= E \left[\mathbf{x}^T (X^T X)^{-1} X^T \boldsymbol{\varepsilon} \cdot \left(X (X^T X)^{-1} \mathbf{x} \right)^T \boldsymbol{\varepsilon} \right] \\ &= E \left[\mathbf{x}^T (X^T X)^{-1} X^T \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T \left(X (X^T X)^{-1} \mathbf{x} \right) \right] \\ &= \mathbf{x}^T (X^T X)^{-1} X^T \underbrace{E \left[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T \right]}_{=\sigma^2 \mathbb{I}_n} \left(X (X^T X)^{-1} \mathbf{x} \right) \\ &= \sigma^2 \mathbf{x}^T (X^T X)^{-1} \mathbf{x} \end{aligned} \tag{24}$$

MEAN SQUARE ERROR (MSE) OF PREDICTION IN REAL TIME

Since $E[e^2] = \sigma^2$, (22), (23) and (24) entail

$$MSE := E[(Y_{n+1} - \hat{y})^2] = \sigma^2 (\mathbf{x}_{n+1}^T (X^T X)^{-1} \mathbf{x}_{n+1} + 1)$$

By (9) and (3)

$$MSE = \mathbf{x}_{n+1}^T C_{\hat{\beta}} \mathbf{x}_{n+1} + \sigma^2 = \text{Var}[\mathbf{x}_{n+1}^T \hat{\beta}] + \sigma^2.$$

By the Gauss-Markov Theorem, there is no linear unbiased one-step predictor in real time with smaller MSE.

AN EXAMPLE OF PREDICTION IN REAL TIME: QUALITY OF WINE GIVEN WEATHER PREDICTOR DATA IN THE CURRENT YEAR ($n + 1$) BEFORE CRUSHING, EXTRACTION, FERMENTATION E.T.C..

We do ordinary multiple regression with three predicting variables x_1, x_2, x_3 , and Y = quality of wine, observed up to year n with

x_1 = Precipitation during the winter months

x_2 = Average temperature during growing season

x_3 = Precipitation during harvesting season

These are now the variables $\mathbf{x}_{n+1}^T = (1, x_{n+1}, x_{n+2}, x_{n+3})$.

MULTIPLE LINEAR REGRESSION: QUALITY OF WINE

Statistical predictor (SPR) for prediction of the annual quality of Bordeaux wine (i.e., before anyone has tasted it) due to Orley Ashenfelter¹ is

$$\hat{y}_{n+1} = \mathbf{x}_{n+1}^T \hat{\beta}, \quad \text{where } \mathbf{x}_{n+1}^T = (1, x_{n+11}, x_{n+12}, x_{n+13})$$

with

$$\hat{\beta} = (12.145, 0.00117, 0.0614, 0.00386)^T$$

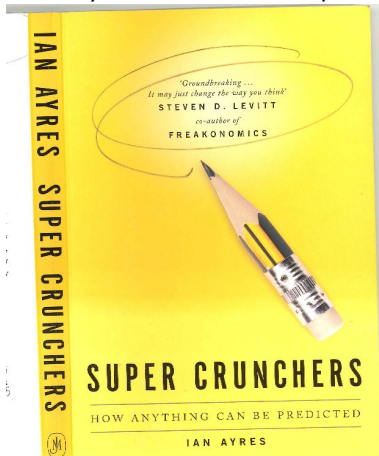
$$\text{Wine Quality} = 12.145 + 0.00117x_1 + 0.0614x_2 - 0.00386x_3$$

This is reported to be a succesful predictor, but is met with resent and embarrasment by many excellent experts on wine tasting.

¹Ashenfelter, Orley: Predicting the quality and prices of Bordeaux wine.
Journal of Wine Economics, 5, 1, 40–52, 2010

A Pocket Book on SPR

For Orley Ashenfelter's SPR for the quality of wine see also I. Ayres: *Super Crunchers. How anything can be predicted*. John Murray (Publishers), Paperback edition 2008, London.



A quote from I.Ayres: Super Crunchers

For a very wide range of prediction problems, statistical prediction rules (SPRs), often rules that are very easy to implement, make predictions that are as reliable as, and typically more reliable than, human experts. The success of SPRs forces us to reconsider our views about what is involved in understanding, explanation, and good reasoning.

(AI ?)

CLIMATE PREDICTION JAPAN METEOROLOGICAL AGENCY (JMA).

For the setting of JMA² we think first of

$$y(t) = \beta_0 + \beta_1 x_1(t) + \cdots + \beta_k x_k(t) + \varepsilon(t)$$

where t is continuous time. $y(t)$ is the temperature ($^{\circ}\text{C}$) in Tokyo at time t . The design: we sample these responses and covariates at n times t_1, \dots, t_n (winter) and set $y_i = y(t_i)$, $x_{ij} = x_j(t_i)$, $\varepsilon_i = \varepsilon(t_i)$ for $i = 1, \dots, n$, $j = 1, \dots, k$. Hence we obtain the ordinary multiple regression equations:

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_k x_{ik} + \varepsilon_i, \quad i = 1, \dots, n,$$

²Source: Statistical Methods for long-range forecast. By Syunji Takahashi / Climate Prediction Division/ JMA

CLIMATE PREDICTION DIVISION JMA

Situation of Multiple Regression Model

	predictand	predictors									
year	Temp	NHEV	FEV	NEZI	FEZI	OKHOTK	MDH	OKINAW	OGASH	WPAH	
1980	-13	-4.0	-9.2	6.1	-24.4	27.9	0.4	13.8	9.6	0.6	
1981	8	-4.9	8.5	-13.6	-4.8	-1.2	0.3	-3.3	-5.0	-2.1	
1982	-21	1.9	13.7	2.1	11.1	-8.2	-11.5	-21.4	-3.8	6.5	
1983	-11	-25.4	-4.8	-12.7	-23.9	-6.1	-11.6	-3.9	2.9	14.9	
1984	7	0.7	14.4	1.0	7.6	-19.9	-5.6	-6.6	-21.9	1.9	
1985	5	16.4	16.6	10.2	5.5	-20.5	-6.8	-12.8	-11.6	-1.7	
1986	-12	-14.8	-12.2	5.6	-28.3	11.1	-10.3	-8.7	-15.2	-2.4	
1987	9	20.2	21.1	-0.4	3.6	-12.2	-2.7	5.9	8.5	0.0	
1988	-19	8.7	14.8	3.2	11.2	33.8	2.6	-2.6	1.3	7.4	
1989	-9	-18.4	-58.7	-17.3	0.8	26.1	2.2	-7.1	-1.7	-0.6	
1990	6	25.1	19.1	10.6	3.5	-4.9	4.0	9.5	1.5	-2.3	
1991	5	5.6	31.4	-7.2	-28.5	-0.8	-0.8	9.6	6.1	-0.9	
1992	0	-20.8	-42.9	3.6	-0.8	-3.9	-4.1	-5.0	3.7	1.9	
1993	-21	30.5	-12.0	-13.9	-24.8	7.1	-12.7	-3.0	2.8	9.3	
1994	24	-33.2	-39.8	14.1	29.4	-14.2	13.8	8.0	-0.5	2.8	
1995	5	-6.4	3.9	-12.8	-3.5	-15.5	-1.0	11.2	11.6	3.9	
1996	6	-32.5	-27.9	4.0	6.2	1.4	5.2	10.6	-3.4	2.8	
1997	4	-3.3	-8.7	-17.4	12.1	-5.3	1.5	-4.5	-2.3	0.4	
1998	5	43.5	33.6	2.1	-27.1	50.0	14.2	14.5	16.0	-5.1	
1999	4	3.2	19.5	14.9	43.6	-1.6	9.2	-6.6	-0.6	-14.3	
2000	15	-3.1	-8.2	-0.1	-2.4	4.1	4.6	-7.4	-4.4	-15.1	

Predictand Vector

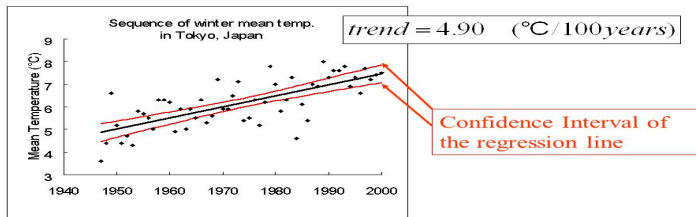
Independent Data

Predictor Matrix

The confidence intervals in the next figure are to be derived later.

SOURCE: STATISTICAL METHODS FOR LONG-RANGE FORECAST. BY SYUNJI TAKAHASHI CLIMATE PREDICTION DIVISION JMA

Property of Calculated Trend



Confidence Interval of the estimated trend

$$3.6 < trend < 6.10 \text{ (}^{\circ}\text{C/100years)}$$

Warming trend in Tokyo is significant

PART III: COEFFICIENT OF DETERMINATION

$$R^2 = \frac{\hat{\beta}^T X^T \mathbf{y} - n\bar{y}^2}{\mathbf{y}^T \mathbf{y} - n\bar{y}^2}$$

R^2 is the fraction of response variance that is captured by the model.

Auxiliaries on LSE Residuals \mathbf{e}_{LSE}

$\varepsilon = \mathbf{Y} - X\beta$ true residuals, unobservable r.v.

$\hat{\varepsilon} = \mathbf{Y} - X\hat{\beta}$ observed LSE residuals as a random vector

$\mathbf{e}_{LSE} = \mathbf{y} - X\hat{\beta} = \mathbf{y} - H\mathbf{y} = \mathbf{y} - \hat{\mathbf{y}}$ observed outcome of $\hat{\varepsilon}$,

$\mathbf{e}_{LSE} = (\hat{\varepsilon}_1, \hat{\varepsilon}_2, \dots, \hat{\varepsilon}_n)^T$.

- $X^T \mathbf{e}_{LSE} = \mathbf{0}_k$ (Check this!). When you look at the scalar product of the first row in X^T and \mathbf{e}_{LSE} this means

$$\sum_{i=1}^n \hat{\varepsilon}_i = 0. \quad (25)$$

- Since $\hat{y}_i = y_i + \hat{\varepsilon}_i$, (25) gives

$$\frac{1}{n} \sum_{i=1}^n \hat{y}_i = \bar{y}. \quad (26)$$

Fundamental Analysis of Variance Identity

LEMMA

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n \hat{e}_i^2 \quad (27)$$

Proof: Let $\mathbf{1}_n$ be the $n \times 1$ vector with all entries equal to 1. We need first to study the $n \times n$ centering matrix C_{ce} defined by

$$C_{ce} := \mathbb{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T.$$

It has been discussed earlier that C_{ce} is idempotent and symmetric.

Fundamental Analysis of Variance Identity

Take next an $n \times 1$ vector \mathbf{a} . Then

$$C_{ce}\mathbf{a} = \mathbb{I}_n\mathbf{a} - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^T\mathbf{a} = \mathbf{a} - \frac{\sum_{i=1}^n a_i}{n} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} a_1 - \bar{a} \\ a_2 - \bar{a} \\ \vdots \\ a_n - \bar{a} \end{pmatrix}. \quad (28)$$

Fundamental Analysis of Variance Identity

We compute the quadratic form $\mathbf{a}^T C_{ce} \mathbf{a}$. Idempotency, symmetry and (28) entail

$$\begin{aligned}\mathbf{a}^T C_{ce} \mathbf{a} &= \mathbf{a}^T C_{ce} C_{ce} \mathbf{a} = \left(C_{ce}^T \mathbf{a} \right)^T C_{ce} \mathbf{a} \\ &= (C_{ce} \mathbf{a})^T C_{ce} \mathbf{a} = \sum_{i=1}^n (a_i - \bar{a})^2,\end{aligned}$$

that is

$$\mathbf{a}^T C_{ce} \mathbf{a} = \sum_{i=1}^n (a_i - \bar{a})^2. \quad (29)$$

Fundamental Analysis of Variance Identity

We have $\mathbf{y} = H\mathbf{y} + \mathbf{e}_{LSE}$, where $\hat{\mathbf{y}} = H\mathbf{y}$. Thus

$$C_{ce}\mathbf{y} = C_{ce}\hat{\mathbf{y}} + C_{ce}\mathbf{e}_{LSE}.$$

Since $\sum_{i=1}^n \hat{e}_i = 0$, see (25) above, we have

$$C_{ce}\mathbf{e}_{LSE} = \mathbf{e}_{LSE} - \frac{1}{n}\mathbf{1}_n \underbrace{\mathbf{1}_n^T \mathbf{e}_{LSE}}_{=\sum_{i=1}^n \hat{e}_i=0} = \mathbf{e}_{LSE}.$$

i.e.,

$$C_{ce}\mathbf{e}_{LSE} = \mathbf{e}_{LSE}. \quad (30)$$

Hence

$$\mathbf{y}^T C_{ce} = \hat{\mathbf{y}}^T C_{ce} + \mathbf{e}_{LSE}^T.$$

Fundamental Analysis of Variance Identity

$\mathbf{y}^T C_{ce} = \hat{\mathbf{y}}^T C_{ce} + \mathbf{e}_{LSE}^T$. Then we multiply and use symmetry, idempotency and (30)

$$\begin{aligned}\mathbf{y}^T C_{ce} C_{ce} \mathbf{y} &= (\hat{\mathbf{y}}^T C_{ce} + \mathbf{e}_{LSE}^T) (C_{ce} \hat{\mathbf{y}} + \mathbf{e}_{LSE}) \\&= \hat{\mathbf{y}}^T C_{ce} C_{ce} \hat{\mathbf{y}} + \hat{\mathbf{y}}^T \underbrace{C_{ce} \mathbf{e}_{LSE}}_{=\mathbf{e}_{LSE}} \\&\quad + \underbrace{\mathbf{e}_{LSE}^T C_{ce} \hat{\mathbf{y}}}_{=\mathbf{e}_{LSE}^T \hat{\mathbf{y}}} + \mathbf{e}_{LSE}^T \mathbf{e}_{LSE} \\&= \hat{\mathbf{y}}^T C_{ce} \hat{\mathbf{y}} + \hat{\mathbf{y}}^T \mathbf{e}_{LSE} + \mathbf{e}_{LSE}^T \hat{\mathbf{y}} + \mathbf{e}_{LSE}^T \mathbf{e}_{LSE}.\end{aligned}\tag{31}$$

Fundamental Analysis of Variance Identity : THE FINAL RESULT

We have found:

$$\mathbf{y}^T C_{ce} C_{ce} \mathbf{y} = \hat{\mathbf{y}}^T C_{ce} \hat{\mathbf{y}} + \hat{\mathbf{y}}^T \mathbf{e}_{LSE} + \mathbf{e}_{LSE}^T \hat{\mathbf{y}} + \mathbf{e}_{LSE}^T \mathbf{e}_{LSE}. \quad (32)$$

- $\mathbf{y}^T C_{ce} C_{ce} \mathbf{y} = \mathbf{y}^T C_{ce} \mathbf{y} = \sum_{i=1}^n (y_i - \bar{y})^2$ by idempotency and (29).
- $\hat{\mathbf{y}}^T C_{ce} \hat{\mathbf{y}} = \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2$, again by (29) and by (26).
- $\mathbf{e}_{LSE}^T \hat{\mathbf{y}} = \hat{\mathbf{y}}^T \mathbf{e}_{LSE} = 0$, since the LSE residuals are orthogonal to $\hat{\mathbf{y}} = H\mathbf{y}$, as found in Lecture 3.
- $\mathbf{e}_{LSE}^T \mathbf{e}_{LSE} = \sum_{i=1}^n \hat{e}_i^2$ by definition of the scalar product.

Hence (27) holds, as claimed. □

Fundamental Analysis of Variance Identity : THE QUADRATIC FORMS

The Fundamental Analysis of Variance Identity in (32) is thus also written as

$$\underbrace{\mathbf{y}^T \mathbf{C}_{ce} \mathbf{y}}_{=SS_T} = \underbrace{\sum_{i=1}^n \left(\hat{y}_i - \bar{\bar{y}} \right)^2}_{SS_R} + \underbrace{\mathbf{e}_{LSE}^T \mathbf{e}_{LSE}}_{=SS_{Res}}. \quad (33)$$

We have found that

$$\sum_{i=1}^n (\hat{y}_i - \bar{\bar{y}})^2 = \sum_{i=1}^n \left(\hat{y}_i - \bar{\bar{y}} \right)^2$$

The decomposition (33) will turn out to be important in Lecture 5, once we can represent SS_R as a quadratic form, too.

Fundamental Analysis of Variance Identity : ILLUSTRATION

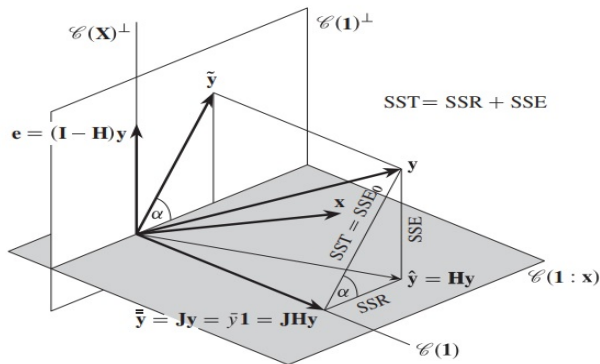


Figure 2 Illustration of $SST = SSR + SSE$.

By Courtesy of: Puntanen, S. and Isotalo, J. and Styan, GPH
Formulas Useful for Linear Regression Analysis and Related Matrix

COEFFICIENT OF DETERMINATION, R^2

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n \hat{e}_i^2$$

The **coefficient of determination** R^2 is defined by

$$R^2 \stackrel{\text{def}}{=} \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

This was defined for simple linear regression in Lecture 1, but we have now seen that this makes sense for \hat{y} computed in multiple regression, too. We shall now show that

$$R^2 = \frac{\hat{\beta}^T X^T \mathbf{y} - n\bar{y}^2}{\mathbf{y}^T \mathbf{y} - n\bar{y}^2} \quad (34)$$

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

The equality $\sum_{i=1}^n (y_i - \bar{y})^2 = \mathbf{y}^T \mathbf{y} - n\bar{y}^2$ is rule (9) in Appendix C of the slides for Lecture 1. By the same rule we get

$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \hat{\mathbf{y}}^T \hat{\mathbf{y}} - n\bar{y}^2.$$

We have $\hat{\mathbf{y}} = H\mathbf{y} = X\hat{\boldsymbol{\beta}}$. Hence

$$\begin{aligned} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 &= \hat{\boldsymbol{\beta}}^T X^T X \hat{\boldsymbol{\beta}} - n\bar{y}^2 \\ &= \hat{\boldsymbol{\beta}}^T X^T X (X^T X)^{-1} X^T \mathbf{y} - n\bar{y}^2 \\ &= \hat{\boldsymbol{\beta}}^T X^T \mathbf{y} - n\bar{y}^2. \end{aligned}$$

Hence we have (34).

We can hence write the regression or model sum of squares as

$$SS_R = \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{y} - n\bar{y}^2. \quad (35)$$

PART IV: STATISTICAL PROPERTIES OF THE LSE RESIDUALS

$\varepsilon = \mathbf{Y} - X\beta_*$ true residuals, unobservable r.v.

$\hat{\varepsilon} = \mathbf{Y} - X\hat{\beta}$ observed LSE residuals as a random vector

The statistical properties of $\hat{\varepsilon}$ are studied in this part of the lecture.

PART IV: STATISTICAL PROPERTIES OF THE LSE RESIDUALS: EXPECTATION

PROPOSITION

$$E[\hat{\epsilon}] = \mathbf{0}_n. \quad (36)$$

Proof: We use the hat matrix to write $\hat{\epsilon} = \mathbf{Y} - X\hat{\beta} = \mathbf{Y} - H\mathbf{Y} = (\mathbb{I}_n - H)\mathbf{Y}$, that is

$$\hat{\epsilon} = (\mathbb{I}_n - H)\mathbf{Y}. \quad (37)$$

Then

$$E[\hat{\epsilon}] = (\mathbb{I}_n - H)E[\mathbf{Y}] =$$

and the rule (1) with the true model

$$= (\mathbb{I}_n - H)(X\beta_* + E[\epsilon])$$

$$= (\mathbb{I}_n - H)X\beta_* = X\beta_* - HX\beta_*$$

But $HX = X(X^T X)^{-1}X^T X = X$. Hence (39) follows.

PART IV: STATISTICAL PROPERTIES OF THE LSE RESIDUALS: COVARIANCE MATRIX

PROPOSITION

$$C_{\hat{\epsilon}} = \sigma^2 (\mathbb{I}_n - H). \quad (38)$$

Proof: By (37) and the rule (2) for covariance matrices of linearly mapped random vectors

$$C_{\hat{\epsilon}} = (\mathbb{I}_n - H) C_Y (\mathbb{I}_n - H)^T$$

But in the true model, $C_Y = \sigma^2 \mathbb{I}_n$. Hence

$$C_{\hat{\epsilon}} = \sigma^2 (\mathbb{I}_n - H) (\mathbb{I}_n - H)^T = \sigma^2 (\mathbb{I}_n - H) (\mathbb{I}_n - H) = \sigma^2 (\mathbb{I}_n - H),$$

where we used the symmetry and idempotency of $\mathbb{I}_n - H$ □

PART IV: STATISTICAL PROPERTIES OF THE LSE RESIDUALS: MULTIVARIATE NORMAL DISTRIBUTION

PROPOSITION

In the normal true model

$$\hat{\varepsilon} \sim N_n \left(\mathbf{0}_n, \sigma^2 (\mathbb{I}_n - H) \right). \quad (39)$$

Proof: The two preceding propositions give the mean vector and covariance matrix as stated. In the true normal model $Y \sim N_n (X\beta_*, \sigma^2 \mathbb{I}_n)$. Since $\hat{\varepsilon} = (\mathbb{I}_n - H) Y$, the random vector $\hat{\varepsilon}$ has a normal distribution. □

REMARK

Since

$$\hat{\varepsilon} \sim N_n \left(\mathbf{0}_n, \sigma^2 (\mathbb{I}_n - H) \right),$$

we find on the main diagonal that

$$\text{Var}(\hat{\varepsilon}_i) = \sigma^2 (1 - h_{ii}).$$

Because a variance is nonnegative, it must hold that $h_{ii} \leq 1$. It will be shown later that $0 < h_{ii} < 1$, since H is idempotent.

DISTRIBUTIONS OF QUADRATIC FORMS OF NORMAL VECTORS

We shall next study the statistical properties of the quadratic form

$$\hat{\varepsilon}^T \hat{\varepsilon}$$

It is shown below that $\hat{\varepsilon} = (\mathbb{I}_n - H) Y = (\mathbb{I}_n - H) \varepsilon$. We choose now to continue with

$$\hat{\varepsilon} = (\mathbb{I}_n - H) \varepsilon,$$

since ε is a normal vector with n independent components. First we collect some facts about $\mathbb{I}_n - H$.

Properties of $\mathbb{I}_n - H$

- $\mathbb{I}_n - H$ is symmetric and idempotent. (Check !) Hence $\mathbb{I}_n - H$ is singular (has no inverse matrix) by an Appendix to Lecture 3.

- Set $A = X(X^T X)^{-1}$, $B = X^T$. By the rule 2. in Appendix B,

$$\text{Tr } H = \text{Tr } AB = \text{Tr } BA = \text{Tr } X^T X (X^T X)^{-1} = \text{Tr } \mathbb{I}_{k+1} = k + 1.$$

Hence by rule 3. in Appendix B,

$$\text{Tr } (\mathbb{I}_n - H) = \text{Tr } \mathbb{I}_n - \text{Tr } H = n - (k + 1). \quad (40)$$

- $\mathbb{I}_n - H$ is positive semidefinite by Appendix XXX below
- Since $\mathbb{I}_n - H$ is positive semidefinite, it follows from (40) by a result in Appendix XXX below that

$$\text{rank } (\mathbb{I}_n - H) = n - (k + 1). \quad (41)$$

DISTRIBUTIONS OF QUADRATIC FORM OF NORMAL VECTORS

- $\mathbf{X} \sim N_n(\mu, \Sigma)$ is an $n \times 1$ Gaussian vector, where Σ is positive definite. The quadratic form is

$$(\mathbf{X} - \mu)^T \Sigma^{-1} (\mathbf{X} - \mu).$$

We use the factorization of Σ^{-1} in (46) of Appendix XXXX to get

$$(\mathbf{X} - \mu)^T \Sigma^{-1} (\mathbf{X} - \mu) = \left(\Sigma^{-1/2} (\mathbf{X} - \mu) \right)^T \Sigma^{-1/2} (\mathbf{X} - \mu).$$

Let $\mathbf{Z} := \Sigma^{-1/2} (\mathbf{X} - \mu)$. Then our rules of computation give that \mathbf{Z} has the mean vector $E[\mathbf{Z}] = \mathbf{0}_n$ and \mathbf{Z} has the covariance matrix

$$\mathbf{C}_Z = \Sigma^{-1/2} \Sigma \Sigma^{-1/2} = \Sigma^{-1/2} \Sigma^{1/2} \Sigma^{1/2} \Sigma^{-1/2} = \mathbf{I}_n.$$

DISTRIBUTIONS OF QUADRATIC FORM OF NORMAL VECTORS

I.e., $\mathbf{Z} \sim N_n(\mathbf{0}_n, \mathbb{I}_n)$. Hence

$$(\mathbf{X} - \mu)^T \Sigma^{-1} (\mathbf{X} - \mu) = \mathbf{Z}^T \mathbf{Z} = \sum_{i=1}^n Z_i^2 \sim \chi^2(n).$$

This is **Thm 9.1.** in Gut, Allan: *An Intermediate Course in Probability. Second Edition*. Note that $z_i \sim N(0, 1)$ are independent and that a finite sum of squares of independent standard normal r.v.'s has $\chi^2(n)$ (chi-squared distribution with n degrees of freedom) (as follows by moment generating functions, c.f., the textbook by Allan Gut).

CHI-SQUARE

DEFINITION

X_1, \dots, X_n are i.i.d., $X_i \sim N(0, 1)$.

$$W = \sum_{i=1}^n X_i^2.$$

W has the **chi-square distribution with n degrees of freedom**, symbolically $W \sim \chi^2(n)$

The pdf of W is

$$f(x; n) = \begin{cases} \frac{x^{\frac{n}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{n}{2}} \Gamma\left(\frac{n}{2}\right)}, & x > 0; \\ 0, & \text{otherwise.} \end{cases}$$

Here $\Gamma(\cdot)$ is the Gamma function.

DISTRIBUTIONS OF A QUADRATIC FORM OF OLS RESIDUALS IN THE NORMAL MULTIPLE REGRESSION

The purpose of all this is to determine the statistical distribution of the unbiased estimator of σ^2 found earlier as the random variable³

$$\widehat{\sigma^2} = \frac{1}{(n - k - 1)} \widehat{\varepsilon}^T \widehat{\varepsilon}.$$

³ $\widehat{\sigma^2}$ in (12) is the outcome of the current quadratic form 

DISTRIBUTION OF THE ESTIMATOR OF VARIANCE IN THE NORMAL MULTIPLE REGRESSION

$$(n - k - 1) \frac{\widehat{\sigma^2}}{\sigma^2} \sim \chi^2(n - k - 1) \quad (42)$$

$\chi^2(n - k - 1)$ is the chi-squared distribution with $n - k - 1$ degrees of freedom.

We shall now establish this important fact.

We have by definition of $\widehat{\sigma^2}$

$$(n - k - 1) \frac{\widehat{\sigma^2}}{\sigma^2} = \frac{\widehat{\epsilon}^T \widehat{\epsilon}}{\sigma^2}.$$

DEGREES OF FREEDOM ?

$$(n - k - 1) \frac{\widehat{\sigma^2}}{\sigma^2} = \frac{\widehat{\boldsymbol{\varepsilon}}^T \widehat{\boldsymbol{\varepsilon}}}{\sigma^2} \sim \chi^2(n - k - 1) \quad (43)$$

is now our claim to be proved. We start by re-checking from Lecture 3 that $\widehat{\boldsymbol{\varepsilon}} = (\mathbb{I}_n - H) \boldsymbol{\varepsilon}$.

By the true model and the hat matrix $H = X(X^T X)^{-1} X^T$

$$\begin{aligned}\hat{\epsilon} &= (\mathbb{I}_n - H) \mathbf{Y} = (\mathbb{I}_n - H) (X\beta_* + \epsilon) \\&= X\beta_* + \epsilon - HX\beta_* - H\epsilon = X\beta_* + \epsilon - X(X^T X)^{-1} X^T X\beta_* - H\epsilon \\&= X\beta_* + \epsilon - X\beta_* - H\epsilon \\&= (\mathbb{I}_n - H) \epsilon\end{aligned}$$

DISTRIBUTIONS OF A QUADRATIC FORM OF OLS RESIDUALS IN THE NORMAL MULTIPLE REGRESSION

Hence

$$(n - k - 1) \frac{\widehat{\sigma^2}}{\sigma^2} = \frac{\widehat{\boldsymbol{\varepsilon}}^T \widehat{\boldsymbol{\varepsilon}}}{\sigma^2} = \frac{\boldsymbol{\varepsilon}^T (\mathbb{I}_n - H)^T (\mathbb{I}_n - H) \boldsymbol{\varepsilon}}{\sigma^2}.$$

The preceding argument about χ -square distribution of quadratic forms has to be revised, when dealing with

$$\frac{\boldsymbol{\varepsilon}^T (\mathbb{I}_n - H)^T (\mathbb{I}_n - H) \boldsymbol{\varepsilon}}{\sigma^2}$$

for the obvious reason that $\mathbb{I}_n - H$ is not invertible.

We quote

PROPOSITION

If $\mathbf{Z} \sim N_n(\mathbf{0}_n, \mathbb{I}_n)$, then

$$\mathbf{Z}^T \mathbf{A} \mathbf{Z} \sim \chi^2(r) \quad (44)$$

if and only if \mathbf{A} is an idempotent matrix with $\text{rank } \mathbf{A} = r$.

This is Corollary 1 to Theorem 5.5 on pp. 117–118 in Rencher, Alvin C and Schaalje, G Bruce: *Linear Models in Statistics*, 2008. The proof in loc.cit. is based on the moment generating function of the quadratic form. Details are omitted here.

DISTRIBUTION OF A QUADRATIC FORM OF OLS RESIDUALS IN THE NORMAL MULTIPLE REGRESSION

We apply this with $A = \mathbb{I}_n - H$ and $\mathbf{Z} := \frac{\boldsymbol{\varepsilon}}{\sigma}$. Then (check this) $\mathbf{Z} \sim N_n(\mathbf{0}_n, \mathbb{I}_n)$ and

$$\frac{\boldsymbol{\varepsilon}^T (\mathbb{I}_n - H)^T (\mathbb{I}_n - H) \boldsymbol{\varepsilon}}{\sigma^2} = \mathbf{Z}^T (\mathbb{I}_n - H)^T (\mathbb{I}_n - H) \mathbf{Z}.$$

By the proposition 7 above we have by idempotency and the previous computation of $\text{rank}(\mathbb{I}_n - H)$

$$\frac{1}{\sigma^2} \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}} = \mathbf{Z}^T (\mathbb{I}_n - H)^T (\mathbb{I}_n - H) \mathbf{Z} = \mathbf{Z}^T (\mathbb{I}_n - H) \mathbf{Z} \sim \chi^2(n - k - 1)$$

as was to be proved. □

APPENDIX A

APPENDIX B : TRACE OF A SQUARE MATRIX

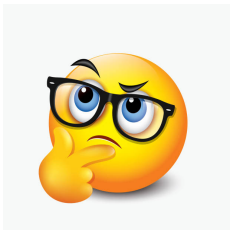
Let A be a square matrix. The **trace** $\text{Tr } A$ of A is the sum of the entries in main diagonal:

$$\text{Tr} \begin{pmatrix} a_{11} & \cdots & a_{1k} \\ \vdots & \ddots & \vdots \\ a_{k1} & \cdots & a_{kk} \end{pmatrix} = \sum_1^k a_{jj}$$

The following facts are easily established; the proofs are left as exercises:

- 1. If A is a $k \times n$ -matrix, and B an $n \times k$ -matrix, then $\text{Tr}(AB) = \text{Tr}(BA)$
- 2. In particular, if a is a column-vector, then $a^T a = \text{Tr}(aa^T)$.
- 3. For any real numbers a and b , $\text{Tr}(aC + bD) = a\text{Tr } C + b\text{Tr } D$

APPENDIX : WONDERINGS



A WONDERING

We have also the identity $\mathbf{Y} = \hat{\mathbf{Y}} + \hat{\boldsymbol{\varepsilon}}$, where $\hat{\mathbf{Y}} = H\mathbf{Y}$. Then

$$\begin{aligned}\hat{\boldsymbol{\varepsilon}} &= (\mathbb{I}_n - H)\mathbf{Y} = \\&= \hat{\mathbf{Y}} + \hat{\boldsymbol{\varepsilon}} - H\hat{\mathbf{Y}} - H\hat{\boldsymbol{\varepsilon}} = \hat{\mathbf{Y}} + \hat{\boldsymbol{\varepsilon}} - HH\mathbf{Y} - H\hat{\boldsymbol{\varepsilon}} \\&= H\mathbf{Y} + \hat{\boldsymbol{\varepsilon}} - H\mathbf{Y} - H\hat{\boldsymbol{\varepsilon}} \\&= (\mathbb{I}_n - H)\hat{\boldsymbol{\varepsilon}} = \hat{\boldsymbol{\varepsilon}},\end{aligned}$$

as $H\hat{\boldsymbol{\varepsilon}} = X(X^T X)^{-1}X^T \hat{\boldsymbol{\varepsilon}} = \mathbf{0}_n$, since

$$\begin{aligned}X^T \hat{\boldsymbol{\varepsilon}} &= X^T (\mathbf{Y} - \hat{\mathbf{Y}}) = X^T (\mathbf{Y} - H\mathbf{Y}) \\&= X^T \mathbf{Y} - X^T H\mathbf{Y} = X^T \mathbf{Y} - X^T X(X^T X)^{-1}X^T \mathbf{Y} = X^T \mathbf{Y} - X^T \mathbf{Y} = \mathbf{0}_{k+1}.\end{aligned}$$

FACTORIZATION AND SQUARE ROOT OF COVARIANCE MATRICES

If Σ is an $n \times n$ symmetric matrix, then Σ can be written as

$$\Sigma = ADA^T,$$

where A is an **orthogonal matrix** ($A^T A = AA^T = I_n$) and D is an $n \times n$ diagonal matrix, with the eigenvalues on the main diagonal. If Σ is a covariance matrix, its eigenvalues λ_i are non-negative. Then

$$\Sigma^{1/2} = AD^{1/2}A^T,$$

where $D^{1/2}$ is an $n \times n$ diagonal matrix, with $\sqrt{\lambda_i}$ on the main diagonal. One checks now that $\Sigma^{1/2}\Sigma^{1/2} = \Sigma$.

FACTORIZATION AND SQUARE ROOT OF COVARIANCE MATRICES

When Σ is positive definite, its eigenvalues are positive and we define

$$\Sigma^{-1/2} = AD^{-1/2}A^T, \quad (45)$$

where $D^{-1/2}$ is an $n \times n$ diagonal matrix, with $1/\sqrt{\lambda_i}$ on the main diagonal. $\Sigma^{-1/2}$ is symmetric, since $D^{-1/2}$ is symmetric. Clearly, $D^{-1/2}D^{-1/2} = D^{-1}$. Then

$$\Sigma^{-1} = \Sigma^{-1/2}\Sigma^{-1/2}. \quad (46)$$

APPENDIX D: ON SYMMETRIC IDEMPOTENT MATRICES

- For A positive semidefinite,

$$\text{rank } A = \text{the number of positive eigenvalues of } A. \quad (47)$$

- For any square A

$$\text{Tr } A = \text{the sum of eigenvalues of } A. \quad (48)$$

If A is a singular, symmetric and idempotent, then A is positive semidefinite.

- *Proof:* By symmetry $A = A^T$, and by idempotency $A^2 = A$.
Then

$$A = A^2 = AA = AA^T.$$

But then

$$\mathbf{x}^T A \mathbf{x} = \mathbf{x}^T A A^T \mathbf{x} = \left(A^T \mathbf{x} \right)^T A^T \mathbf{x} = \| A \mathbf{x} \|^2 \geq 0.$$



EIGENVALUES OF A SYMMETRIC AND IDEMPOTENT MATRIX

If A is an $n \times n$ symmetric and idempotent matrix with $\text{rank } A = r$, then A has r eigenvalues equal to 1 and $n - r$ eigenvalues equal to 0.

- *Proof:* Let \mathbf{x} satisfy $A\mathbf{x} = \lambda\mathbf{x}$. Then

$$A^2\mathbf{x} = A(A\mathbf{x}) = \lambda A\mathbf{x} = \lambda^2\mathbf{x}.$$

Also

$$A^2\mathbf{x} = A\mathbf{x} = \lambda\mathbf{x}.$$

That is, $\lambda^2\mathbf{x} = \lambda\mathbf{x} \Leftrightarrow (\lambda - \lambda^2)\mathbf{x} = \mathbf{0}$. But an eigenvector is not the zero vector. Hence $(\lambda - \lambda^2) = \lambda(1 - \lambda) = 0$, which holds for $\lambda = 0$ and $\lambda = 1$.

Since A is positive semidefinite, by (47) there are r eigenvalues equal to 1 and $n - r$ eigenvalues equal to 0.



EIGENVALUES OF A SYMMETRIC AND IDEMPOTENT MATRIX

If A is an $n \times n$ symmetric and idempotent matrix with $\text{rank } A = r$, then $\text{Tr } A = r$.

- *Proof:* This follows by (48), as by the preceding statement the sum eigenvalues of A is r . □

The following can be established without direct reference to momentgenerating functions, as shown next.

PROPOSITION

If $\mathbf{Z} \sim N_n(\mathbf{0}_n, \mathbb{I}_n)$ and A is an idempotent matrix with $\text{rank } A = r$, then

$$\mathbf{Z}^T A \mathbf{Z} \sim \chi^2(r). \quad (49)$$

Proof: Söderström -Stoica: System Identification. Prentice Hall, 1986.