

INTRODUCTION

REGRESSION MODEL IS A **STATISTICAL METHOD** FOR PREDICTING VALUES OF ONE OR MORE **RESPONSE VARIABLES** FROM A SET OF **PREDICTOR VARIABLES** ON THE RESPONSES.

MATHEMATICAL FORMULATION

WITH r PREDICTOR VARS
 n RESPONSE VARS

CLASSICAL LINEAR REG. MODEL

$$\hat{Y} = Z\beta + \varepsilon$$

IN WHICH, $E(\varepsilon) = 0$ AND $\text{cov}(\varepsilon) = \sigma^2 I$.

OR MORE EXPLICITY, WE CAN WRITE THE MODEL AS :

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 & z_{11} & z_{12} & \dots & z_{1r} \\ 1 & z_{21} & z_{22} & \dots & z_{2r} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & z_{n2} & \dots & z_{nr} \end{bmatrix}_{n \times (r+1)} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_r \end{bmatrix}_{(r+1) \times 1} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}_{n \times 1}$$

REMARK AT THIS STAGE WE DON'T ASSUME THAT ε IS DISTRIBUTED AS NORMAL; WE ONLY NEED IT TO BE ZERO-MEAN AND COMPONENT-WISELY UNCORRELATED.

METHODS FOR DERIVING ESTIMATION

WE WANT TO PREDICT THE RESPONSE FOR GIVEN VALUES OF PREDICTOR VARIABLES, THEN IT HAS TO "FIT" THE MODEL DESCRIBED AS ABOVE. WE MAINLY HAVE TWO IDEAS TO REACH OUR GOAL :

1) LEAST SQUARES ESTIMATION

IDEA: MINIMIZE THE DISTANCE BETWEEN THE OBSERVED RESPONSES AND PREDICTIONS OF RESPONSES.

OPTIMIZATION MODEL :

OPTIMIZATION MODEL:

$$(L.E.S.) \quad \hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{r+1}} (y - z\beta)^T (y - z\beta)$$

RESULT 1 THE SOLUTION OF (L.E.S.) IS $\hat{\beta} = (z^T z)^{-1} z^T y$

PROOF JUST BY BASIC COMPUTATION:

$$\begin{aligned} f(\beta) &:= (y - z\beta)^T (y - z\beta) \\ &= (y^T - \beta^T z^T)(y - z\beta) \\ &= y^T y - y^T z\beta - \beta^T z^T y + \beta^T z^T z\beta \\ \text{AND } (y^T z\beta)^T &= \beta^T z^T y = y^T y - 2\beta^T z^T y + \beta^T z^T z\beta \end{aligned}$$

MOREOVER, WE STATE THAT:

$$\textcircled{1} \quad \frac{\partial \beta^T z^T y}{\partial \beta} = z^T y \quad \text{AND} \quad \textcircled{2} \quad \frac{\partial \beta^T z^T z\beta}{\partial \beta} = 2z^T z\beta$$

TILL NOW THE DERIVATING WITH RESPECT TO A VECTOR IS NOT RIGOROUSLY DEFINED, THUS WE HAVE TO SHOW THE DETAIL

$$\begin{aligned} \text{SUB-PROOF OF } \textcircled{1} \quad \beta^T z^T y &= (\beta_1, \dots, \beta_{r+1}) \begin{pmatrix} (z^T y)_1 \\ \vdots \\ (z^T y)_{r+1} \end{pmatrix} \\ &= \sum_{i=1}^{r+1} \beta_i (z^T y)_i \\ \Rightarrow \frac{\partial \beta^T z^T y}{\partial \beta_i} &= (z^T y)_i \end{aligned}$$

\Rightarrow WE HAVE THE COMPATIBLE NOTATION:

$$\frac{\partial \beta^T z^T y}{\partial \beta} := z^T y$$

SUB-PROOF OF $\textcircled{2}$ OBVIOUSLY $z^T z$ IS A SYMMETRIC MATRIX,
WE SHALL DENOTE $A := z^T z$.

$$\Rightarrow \beta^T A \beta = \sum_{i=1}^{r+1} \beta_i A \beta_i \beta_i$$

WE SHALL DENOTE $A := ZZ'$.

$$\Rightarrow \underline{\beta}^T A \underline{\beta} = \sum_{i,j=1}^{r+1} \beta_i A_{ij} \beta_j$$

$$\Rightarrow \frac{\partial \underline{\beta}^T A \underline{\beta}}{\partial \beta_i} = 2 \cdot \sum_{j=1}^{r+1} A_{ij} \beta_j = 2 (A\underline{\beta})_i$$

\Rightarrow IN THIS CASE WE ALSO HAVE THE COMPATIBLE NOTATION:

$$\frac{\partial \underline{\beta}^T A \underline{\beta}}{\partial \underline{\beta}} := 2 \underline{Z}^T \underline{Z} \underline{\beta}$$

RETURN BACK TO $f(\underline{\beta}) = \underline{y}^T \underline{y} - 2 \underline{\beta}^T \underline{Z}^T \underline{y} + \underline{\beta}^T \underline{Z}^T \underline{Z} \underline{\beta}$, AND LET

$$\frac{\partial f}{\partial \underline{\beta}} = -2 \underline{Z}^T \underline{y} + 2 \underline{Z}^T \underline{Z} \underline{\beta} = 0$$

$$\Rightarrow \underline{\beta} = (\underline{Z}^T \underline{Z})^{-1} \underline{Z}^T \underline{y} =: \hat{\underline{\beta}}$$

BESIDES PROVING $\hat{\underline{\beta}}$ IS A STATIONARY POINT OF f , WE HAVE ALSO TO STUDY THE SECOND ORDER DERIVATIVE:

$$\frac{\partial^2 f}{\partial \underline{\beta}^2} = 2 \underline{Z}^T \underline{Z}$$

RECALL A VERY IMPORTANT LEMMA WITHOUT PROOF:

WITH THE CLASSICAL DEF: POSITIVE DEFINITIVE = SYM. + $\forall x > 0$,

LEMMA 2 A REAL SYMMETRIC MATRIX A IS POSITIVE DEFINITIVE IF AND ONLY IF THERE EXISTS A

NON-SINGULAR MATRIX Z SUCH THAT

$$A = Z^T Z.$$

THANKS TO THIS LEMMA; AND WITH THE THEORY IN BASIC OPTIMIZATION, WE KNOW: FOR OBTAINING THE MINIMUM AT $\hat{\underline{\beta}}$, IT REQUIRES $Z^T Z$ TO BE POSITIVE DEFINITIVE, WHICH IS EQUIVALENT TO REQUIRING Z TO BE FULL-RANK.

2) MAXIMUM LIKELIHOOD ESTIMATION

IDEA: CHOOSE THE PARAMETER WHICH MAXIMIZE THE LIKELIHOOD AS THE ESTIMATE OF YOUR PARAMETER

IDEA: CHOOSE THE PARAMETER WHICH MAXIMIZE THE LIKELIHOOD AS THE ESTIMATION OF TRUE PARAMETER.

REMARK FOR STUDYING THE LIKELIHOOD, WE MUST INTRODUCE THE PROBABILISTIC DISTRIBUTIONS INTO THE MODEL. THEREFORE, FROM NOW ON WE ASSUME

$$\varepsilon \sim N(\Omega, \sigma^2 I) \quad \leftarrow \text{N - multi-variate Gaussian distribution}$$

THE MODEL IS AS $y = Z\beta + \varepsilon$, WHICH IMPLIES $y \sim N(Z\beta, \sigma^2 I)$, FROM WHICH WE CAN CONSTRUCT THE MAXIMUM LIKELIHOOD FUNCTION

$$\begin{aligned} L(y|\beta) &= \frac{1}{(2\pi)^{n/2} |\sigma^2 I|^{1/2}} e^{-\frac{(y-Z\beta)'(y-Z\beta)}{2\sigma^2}} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot e^{-\frac{(y-Z\beta)'(y-Z\beta)}{2\sigma^2}} \end{aligned}$$

THEN LET $\hat{\beta} = \arg \max_{\beta \in \mathbb{R}^{r+1}} L(y|\beta)$, WHICH IS EXACTLY EQUIVALENT TO

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{r+1}} (y-Z\beta)'(y-Z\beta)$$

THEREFORE WE WILL GET THE SAME RESULT AS IN THE LEAST SQUARES ESTIMATION.

PROPERTIES OF ESTIMATORS

LAST SECTION WE HAVE USED TWO METHODS, OR WE SHALL SAY, TWO DIFFERENT PHILOSOPHIES TO DERIVE THE ESTIMATOR

$$\hat{\beta} = (Z'Z)^{-1} Z'y$$

HERE WE STILL JUST ASSUME $E[\varepsilon] = \Omega$ AND $\text{cov}(\varepsilon) = \sigma^2 I$, RATHER THAN GIVE THE ERROR ANY SPECIFIC DISTRIBUTION.

MOREOVER, HERE WE INTRODUCE THE RESIDUAL $\hat{\varepsilon}$

MOREOVER, HERE WE INTRODUCE THE RESIDUAL $\hat{\varepsilon}$ AS THE ESTIMATOR OF ERROR ε . IT IS DEFINED AS:

$$\hat{\varepsilon} := y - \hat{y} = y - z(z'z)^{-1}z'y = [I - z(z'z)^{-1}z']y$$

THEN WE HAVE THE FOLLOWING RESULT:

RESULT 3 UNDER THE BASIC ASSUMPTION ON ε , THE ESTIMATOR $\hat{\beta} = (z'z)^{-1}z'y$ HAS

$$E[\hat{\beta}] = \beta \quad \text{UNBIASED ESTIMATOR}$$

AND

$$\text{cov}(\hat{\beta}) = \sigma^2 (z'z)^{-1}$$

THE RESIDUAL $\hat{\varepsilon}$ HAS

$$E[\hat{\varepsilon}] = 0$$

AND

$$\text{cov}(\hat{\varepsilon}) = \sigma^2 [I - z(z'z)^{-1}z']$$

MOREOVER, $\hat{\beta}$ AND $\hat{\varepsilon}$ ARE UNCORRELATED.

PROOF

$$\begin{aligned} E[\hat{\beta}] &= E[(z'z)^{-1}z'y] \\ &= E[(z'z)^{-1}z'(z\beta + \varepsilon)] \\ &= E[(z'z)^{-1}z'z\beta] + E[(z'z)^{-1}z'\varepsilon] \\ &= \beta \quad \text{IN OUR SETTING, ALL THE THINGS ARE DETERMINISTIC, EXCEPT } \varepsilon. \end{aligned}$$

$$\begin{aligned} \text{cov}(\hat{\beta}) &= E[(\hat{\beta} - E[\hat{\beta}])(\hat{\beta} - E[\hat{\beta}])'] \\ &= E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)'] \\ &= E[\hat{\beta}\hat{\beta}' - \hat{\beta}\beta' - \beta\hat{\beta}' + \beta\beta'] \\ &= E[(z'z)^{-1}z'y \cdot y'z(z'z)^{-1}] + \beta\beta' \\ &\quad - E[(z'z)^{-1}z'y \cdot \beta' + \beta \cdot y'z(z'z)^{-1}] \end{aligned}$$

RECALL THAT $y = z\beta + \varepsilon$

$$\begin{aligned} &= E[(z'z)^{-1}z'(z\beta + \varepsilon)(z\beta + \varepsilon)'z(z'z)^{-1}] + \beta\beta' \\ &\quad - E[(z'z)^{-1}z'(z\beta + \varepsilon)\beta'] + \beta \cdot (z\beta + \varepsilon)'z(z'z)^{-1} \\ &= E[(z'z)^{-1}(z\beta \cdot \beta' + z\beta \cdot \varepsilon' + \varepsilon \cdot \beta' + \varepsilon\varepsilon')z(z'z)^{-1}] + \beta\beta' \end{aligned}$$

$$\begin{aligned}
 & - E[(\bar{Z}\bar{Z})^{-1} (\bar{Z}\beta + \bar{\varepsilon}) \cdot \beta' + \beta \cdot (\bar{Z}\beta + \bar{\varepsilon}) \bar{Z}(\bar{Z}\bar{Z})^{-1}] \\
 &= E[(\bar{Z}'\bar{Z})^{-1} (\bar{Z}\beta \cdot \beta' \bar{Z} + \bar{Z}\beta \cdot \bar{\varepsilon}' + \bar{\varepsilon}' \cdot \beta' \bar{Z} + \bar{\varepsilon}' \bar{\varepsilon}) \bar{Z}(\bar{Z}'\bar{Z})^{-1} + \beta \beta'] \\
 &\quad - E[(\bar{Z}'\bar{Z})^{-1} \bar{Z}' \bar{Z} \beta \beta' + \beta' \beta \bar{Z}' \bar{Z} (\bar{Z}\bar{Z})^{-1}] \\
 &= E[(\bar{Z}'\bar{Z})^{-1} \bar{Z}' \bar{Z} \beta \cdot \beta' \bar{Z} (\bar{Z}'\bar{Z})^{-1}] + E[(\bar{Z}'\bar{Z})^{-1} \bar{Z}' \bar{Z} \bar{\varepsilon}' (\bar{Z}'\bar{Z})^{-1}] \\
 &\quad + \beta \beta' - \beta \beta' - \beta \beta' \\
 &= (\bar{Z}'\bar{Z})^{-1} \bar{Z}' \underbrace{E[\bar{\varepsilon}' \bar{\varepsilon}]}_{\text{LUCKILY, } \bar{\varepsilon} \text{ HAS ZERO MEAN}} \bar{Z} (\bar{Z}'\bar{Z})^{-1} \\
 &= (\bar{Z}'\bar{Z})^{-1} \bar{Z}' \sigma^2 I \bar{Z} (\bar{Z}'\bar{Z})^{-1} \\
 &= \sigma^2 (\bar{Z}'\bar{Z})^{-1}
 \end{aligned}$$

LUCKILY, ξ
HAS ZERO MEAN,
THIS EXACTLY IS
 ξ 'S COVARIANCE
MATRIX.

NOW WE CAN EASILY OBTAIN THAT

$$\begin{aligned} \mathbb{E}[\hat{\zeta}] &= \mathbb{E}[y - \hat{\alpha}\hat{\beta}] \\ &= \mathbb{E}[y - \bar{z} \cdot (\bar{z}'\bar{z})^{-1}\bar{z}'y] \\ &= \bar{z}'\bar{\beta} - \bar{z} \cdot (\bar{z}'\bar{z})^{-1}\bar{z}'\bar{z}'\bar{\beta} = 0 \end{aligned}$$

ONE CAN ALSO SHOW THAT $\text{Cov}(\hat{\beta}) = \sigma^2 [I - Z(Z'Z)^{-1}Z']$ AND $\text{Cov}(\hat{\beta}, \hat{\varepsilon}) = 0$ BY THE SIMILAR COMPUTATION, BUT WE OMIT THE DETAIL HERE.

GEOOMETRY OF LEAST SQUARES

A GEOMETRICAL INTERPRETATION OF LEAST SQUARES TECHNIQUE HIGHLIGHTS THE NATURE OF THE CONCEPT. FIRST WE SHALL REFRESH US WITH THE NOTATION OF PROJECTION IN THE HILBERT SPACE. IN FINITELY DIMENSIONAL CASE, THE PROJECTION OPERATOR CAN BE REPRESENTED AS CANONICAL MATRIX. WE GIVE THE FOLLOWING DEFINITION:

DEF 4 A REAL-VALUED MATRIX P IS CALLED **ORTHOGONAL PROJECTION MATRIX** IF IT IS:

- 1) SYMMETRIC,
 - 2) $P^2 = P$.

IN OUR CASE, BOTH RESPONSE VARIABLE y AND PREDICTED RESPONSES VARIABLE \hat{y} BELONG TO FINITE DIMENSIONAL HILBERT SPACE H^n . THE MAPPING WHICH IS THE LINEAR

RESPONSES VARIABLE \hat{y} BELONG TO FINITE DIMENSIONAL HILBERT SPACE \mathbb{R}^n , THE LINK BETWEEN THEM IS THE LINEAR OPERATOR (OR WE CALL IT MATRIX):

$$\hat{y} = \mathbf{z}\hat{\beta} = \mathbf{z} \cdot (\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}'y \quad \text{← ALWAYS ASSUME THAT } \mathbf{z} \text{ IS OF FULL RANK}$$

WE DENOTE $P := \mathbf{z}(\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}'$, AND OBVIOUSLY:

$$P: \mathbb{R}^n \rightarrow \mathbb{R}^n, P: y \mapsto \mathbf{z}(\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}'y$$

THEN, WE CAN CHECK THAT

$$P^1 = [\mathbf{z}(\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}']' = \mathbf{z}(\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}' = P$$

$$P^2 = P \cdot P = \mathbf{z}(\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}' \cdot \mathbf{z}(\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}' = \mathbf{z}(\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}' = P$$

ACCORDING TO DEFINITION Φ , WE KNOW THAT $P = \mathbf{z}(\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}'$ IS A PROJECTION OPERATOR (MATRIX) ON \mathbb{R}^n . SINCE WE HAS ASSUMED THAT \mathbf{z} IS OF FULL RANK, RECALL LEMMA 2, WE KNOW THAT $\mathbf{z}'\mathbf{z}$ MUST BE POSITIVE DEFINITE. THEN WE CAN USE THE POWERFUL TOOL: **SPECTRAL THEOREM**, FOR STUDY THE PROPERTY OF OUR PROJECTION MATRIX P .

THEOREM 8 [REAL SPECTRAL THEOREM IN MATRIX PRESENTATION]
SUPPOSE P IS AN $n \times n$ MATRIX, THEN \mathbb{R}^n HAS AN ORTHONORMAL BASIS CONSISTING OF EIGENVECTOR OF P IF AND ONLY IF P IS SYMMETRIC.

↑ FOR THE PROOF PLEASE CHECK THE BOOK OF AXLER - LINEAR ALGEBRA DONE RIGHT.

THEN, WE CAN WRITE

$$\mathbf{z}'\mathbf{z} = \sum_{i=1}^{r+1} \lambda_i e_i e_i'$$

WHERE $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{r+1} > 0$ ARE THE EIGENVALUES OF $\mathbf{z}'\mathbf{z}$ AND e_1, e_2, \dots, e_{r+1} ARE THE CORRESPONDING EIGENVECTORS. FURTHER,

$$(\mathbf{z}'\mathbf{z})^{-1} = \sum_{i=1}^{r+1} \frac{1}{\lambda_i} e_i e_i'$$

CONSIDER $g_i = \lambda_i^{-1/2} \mathbf{z} e_i$, WHICH IS A LINEAR COMBINATION OF THE COLUMNS OF \mathbf{z} . THEN WE HAVE

CONSIDER $q_i = \lambda_i^{-1/2} Z e_i$, WHICH IS A LINEAR COMBINATION OF THE COLUMNS OF Z . THEN WE HAVE

$$\begin{aligned} q_i^T q_k &= \lambda_i^{-1/2} e_i^T Z^T \cdot \lambda_k^{-1/2} Z e_k \\ &= \lambda_i^{-1/2} \lambda_k^{-1/2} e_i^T \left(\sum_{j=1}^{r+1} \lambda_j e_j e_j^T \right) e_k \\ &= \begin{cases} 0 & \text{if } i \neq k \\ 1 & \text{if } i = k \end{cases} \end{aligned}$$

THAT IS, THE $r+1$ VECTORS q_i ARE MUTUALLY PERPENDICULAR AND HAVE THE UNIT LENGTH. THEIR LINEAR COMBINATION SPAN THE SPACE OF ALL LINEAR COMBINATIONS OF THE COLUMNS OF Z . MOREOVER,

$$Z(Z^T Z)^{-1} Z^T = \sum_{i=1}^{r+1} \lambda_i^{-1} Z e_i e_i^T Z^T = \sum_{i=1}^{r+1} q_i q_i^T$$

THE PROJECTION OF y ON A LINEAR COMBINATION OF $\{q_1, \dots, q_{r+1}\}$ IS

$$\sum_{i=1}^{r+1} (q_i^T y) q_i = \left(\sum_{i=1}^{r+1} q_i q_i^T \right) y = Z(Z^T Z)^{-1} Z^T y = Z \hat{y}$$

THUS, THE MULTIPLICATION BY $Z(Z^T Z)^{-1} Z^T$ PROJECTS A VECTOR ONTO THE SPACE SPANNED BY THE COLUMNS OF Z .

EVALUATION: SUM-OF-SQUARES DECOMPOSITION

FIRST WE CAN OBSERVE THAT

$$\begin{aligned} Z \hat{y} &= Z^T (y - \hat{y}) \\ &= Z^T [I - Z(Z^T Z)^{-1} Z^T] y \\ &= [Z^T - Z^T Z(Z^T Z)^{-1} Z^T] y = 0 \end{aligned} \tag{1}$$

AND $\hat{y}^T \hat{y} = \hat{\beta}^T Z^T \hat{y} = 0$.

THEN WE CAN DIRECTLY GET THAT

$$\begin{aligned} y^T y &= (\hat{y} + y - \hat{y})^T (\hat{y} + y - \hat{y}) \\ &= (\hat{y}^T + \hat{\beta}^T)(\hat{y} + \hat{\beta}) \\ &= \hat{y}^T \hat{y} + \hat{y}^T \hat{\beta} + \hat{\beta}^T \hat{y} + \hat{\beta}^T \hat{\beta} \\ &= \hat{y}^T \hat{y} + \hat{\beta}^T \hat{\beta} \end{aligned} \tag{2}$$

$$= \hat{y}_1 \hat{y}_1 + \hat{y}_2 \hat{y}_2 + \dots + \hat{y}_n \hat{y}_n + \underline{\hat{\epsilon} \hat{\epsilon}}$$

$$= \hat{y}_1^2 + \hat{y}_2^2 + \dots + \hat{y}_n^2$$

FURTHERMORE, SINCE THE FIRST COLUMN OF Z IS $\underline{1}$, THEN ① TELLS US ALSO $\underline{1}' \underline{\hat{\epsilon}} = 0$, THEN

$$0 = \underline{1}' \underline{\hat{\epsilon}} = \sum_{j=1}^n \hat{\epsilon}_j = \sum_{j=1}^n (y_j - \hat{y}_j) = \sum_{j=1}^n y_j - \sum_{j=1}^n \hat{y}_j \quad ③$$

$$\Leftrightarrow \bar{y} = \hat{\bar{y}}, \text{ if we define } \bar{y} := \frac{1}{n} \sum_{j=1}^n y_j \text{ AND } \hat{\bar{y}} = \frac{1}{n} \sum_{j=1}^n \hat{y}_j$$

COMBINE ② AND ③ WE CAN GET

$$y'y - n\bar{y}^2 = \hat{y}' \hat{y} - n(\bar{y})^2 + \underline{\hat{\epsilon}}' \underline{\hat{\epsilon}}$$

$$\Leftrightarrow \sum_{j=1}^n (y_j - \bar{y})^2 = \sum_{j=1}^n (\hat{y}_j - \bar{y})^2 + \sum_{j=1}^n \hat{\epsilon}_j^2$$

\Rightarrow TOTAL sum of REGRESSION RESIDUAL
SQUARES ABOUT = sum of SQUARES + sum of MEAN SQUARES

THE PRECEDING SUM OF SQUARES DECOMPOSITION SUGGESTS THAT THE QUALITY OF THE MODELS FIT CAN BE MEASURED BY COEFFICIENT OF DETERMINATION:

$$R^2 = 1 - \frac{\sum_{j=1}^n \hat{\epsilon}_j^2}{\sum_{j=1}^n (y_j - \bar{y})^2} = \frac{\sum_{j=1}^n (\hat{y}_j - \bar{y})^2}{\sum_{j=1}^n (y_j - \bar{y})^2}$$

ONE CAN OBSERVE THAT $R^2 = 1$ IF THE FITTED EQUATION PASSES THROUGH ALL THE DATA POINTS; $R^2 = 0$ IF $\hat{\beta}_0 = \bar{y}$ AND $\hat{\beta}_1 = \dots = \hat{\beta}_r = 0$, THAT IS, THE PREDICTOR VARIABLES IN FACT HAS NO INFLUENCE ON THE RESPONSE VARIABLES.

INFERENCES ABOUT REGRESSION MODEL

FROM NOW ON WE FORMALLY INTRODUCE A SPECIFIC GAUSSIAN DISTRIBUTION TO THE ERROR $\underline{\hat{\epsilon}}$, AS

$$\underline{\hat{\epsilon}} \sim N_n(0, \sigma^2 I)$$

IN THE ROOT, THE "n"
REPRESENTS THE
DIMENSION

THEN WE HAVE THE FOLLOWING IMPORTANT RESULT:

THEN WE HAVE THE FOLLOWING IMPORTANT RESULT:

RESULT 6 LET $\mathbf{y} = \mathbf{z}\beta + \varepsilon$, WHERE \mathbf{z} HAS FULL RANK $r+1$ AND ε IS DISTRIBUTED AS $\varepsilon \sim N_{n \times 1}(\mathbf{0}, \sigma^2 I)$, THEN THE MAXIMAL LIKELIHOOD ESTIMATOR OF β IS AS THE LEAST SQUARES ESTIMATOR $\hat{\beta}$. MOREOVER,

$$1) \hat{\beta} = (\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}'\mathbf{y} \sim N_{r+1}(\beta, \sigma^2(\mathbf{z}'\mathbf{z})^{-1})$$

$$2) \hat{\varepsilon} = \mathbf{y} - \mathbf{z}\hat{\beta}, \hat{\varepsilon} \perp \hat{\beta}. \text{ FURTHER } n\hat{\sigma}^2 = \hat{\varepsilon}'\hat{\varepsilon} \sim \sigma^2 \chi_{n-r-1}^2 \quad \textcircled{1}$$

WHERE $\hat{\sigma}^2$ IS THE MAXIMAL LIKELIHOOD ESTIMATOR OF σ^2 .

PROOF

WE HAVE SHOWN THAT THE ESTIMATOR DERIVED BY LEAST SQUARES ESTIMATION AND MAXIMAL LIKELIHOOD METHOD COINCIDE UNDER THE SPECIFIC ASSUMPTION ON ε .

FOR THE POINT 1, IT IS A DIRECT CONCUSSION OF RESULT 3.
FOR THE POINT 2, ALSO FROM RESULT 3 WE KNOW THAT $\text{cov}(\hat{\varepsilon}, \hat{\beta}) = 0$, AND UNDER THE ASSUMPTION OF GAUSSIAN DISTRIBUTION, NOW THEY ARE ALSO INDEPENDENT. THEREFORE, WE ONLY HAVE TO PROVE STATEMENT ①, THIS NEEDS SOME CALCULATIONS:

$$\mathbf{P} := \mathbf{z}(\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}'$$

$$\hat{\varepsilon} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{P})\mathbf{y}, \text{ WHERE } \hat{\mathbf{y}} \sim N_r(\mathbf{z}\beta, \sigma^2 I)$$

$$\Rightarrow \hat{\varepsilon} \sim N_n([(\mathbf{I} - \mathbf{P})\mathbf{z}\beta], (\mathbf{I} - \mathbf{P})^T \sigma^2 I (\mathbf{I} - \mathbf{P}))$$

$$\text{NOTICE } [(\mathbf{I} - \mathbf{P})\mathbf{z}\beta] = (\mathbf{I} - \mathbf{z}(\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}')\mathbf{z}\beta = \mathbf{z}\beta - \mathbf{z}\beta = \mathbf{0},$$

$$(\mathbf{I} - \mathbf{P})^T \sigma^2 I (\mathbf{I} - \mathbf{P}) = \sigma^2 (\mathbf{I} - \mathbf{P})$$

$$\Rightarrow \hat{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2 (\mathbf{I} - \mathbf{P}))$$

P IS PROJ. MATRIX, THEN $\mathbf{I} - \mathbf{P}$ ALSO IS,
ONE CAN USE DEFINITION TO CHECK.
GEOMETRICALLY SPEAKING, IT'S A
PROJ. UP. TO THE PERPENDICULAR SPACE.

WE DO ALREADY KNOW THAT $\mathbf{I} - \mathbf{P}$ IS ALSO A PROJECTION MATRIX,
NOW WE HAVE TO USE AN IMPORTANT PROPERTY OF EIGENVALUES
AND EIGENVECTORS OF A PROJECTION MATRIX:

SUPPOSE H IS AN (ORTHOGONAL) PROJECT MATRIX, WE SHALL STUDY ITS EIGENVALUES AND EIGENVECTORS AS:

SUPPOSE H IS AN (ORTHOGONAL) PROJECT MATRIX, WE SHALL STUDY ITS EIGENVALUES AND EIGENVECTORS AS:

$$H\underline{x} = \lambda \underline{x} \text{ FOR SOME } \underline{x} \neq 0 \text{ AND } \lambda \in \mathbb{R}$$

SINCE H IS IDEMPOTENT, THAT IS $H = H^2$, THEREFORE:

$$\lambda \underline{x} = H\underline{x} = H^2 \underline{x} = H \cdot H\underline{x} = H \cdot \lambda \underline{x} = \lambda^2 \underline{x} \text{ FOR SOME } \underline{x} \neq 0.$$

$\Rightarrow \lambda$ MUST BE EITHER 0 OR 1.

ALSO THANKS TO REAL SPECTRAL THEOREM, SINCE $I - P$ IS REAL AND SYMMETRIC, WE CAN FIND ORTHONORMAL BASIS OF \mathbb{R}^n CONSISTING OF EIGENVECTORS OF $I - P$. WE ALREADY KNOW THAT THE EIGENVALUES OF $I - P$ ARE EITHER 0 OR 1, ACCORDING TO THIS NICE PROPERTY, WE CAN STUDY THE DIMENSION OF CORRESPONDING EIGENSPACE V_0 AND V_1 ($V_0 \oplus V_1 = \mathbb{R}^n$) BY JUST COMPUTING THE TRACE OF $I - P$:

$$\text{Tr}(I - P) = \text{Tr}(I) - \text{Tr}(P)$$

$$\begin{aligned} & \stackrel{\text{CYCLIC PROPERTY}}{=} \text{Tr}(I) - \text{Tr}(Z(Z^T Z)^{-1} Z) \\ & \stackrel{\text{OF TRACE}}{=} \text{Tr}(I) - \text{Tr}((Z^T Z)^{-1} Z^T Z) \\ & = \text{Tr}(I) - \text{Tr}(\tilde{I}) \end{aligned}$$

NOTICE THAT I AND \tilde{I} ARE BOTH IDENTITY MATRIX (OPERATOR), BUT $I: \mathbb{R}^n \rightarrow \mathbb{R}^n$; $\tilde{I}: \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$, THEREFORE

$$\text{Tr}(I - P) = n - (r + 1) = n - r - 1$$

THIS TELLS US $\dim V_0 = r + 1$, $\dim V_1 = n - r - 1$, WHERE 0 AND 1 ARE CORRESPONDING EIGENVALUES OF MATRIX $I - P$. THEN AGAIN IN TERMS OF REAL SPECTRAL THEOREM, WE HAVE THE FOLLOWING EQUALITY:

$$I - P = Q \Lambda Q'$$

WHERE:

$$\Lambda = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix}, \quad Q = [\underline{e}_1, \underline{e}_2, \dots, \underline{e}_n]$$

$\lambda_1, \lambda_2, \dots, \lambda_n$ ARE EIGENVALUES OF MATRIX $I - P$, IN PARTICULAR

$\lambda_1, \lambda_2, \dots, \lambda_n$

$\lambda_1, \lambda_2, \dots, \lambda_n$ ARE EIGENVALUES OF MATRIX $I - P$, IN PARTICULAR, $r+1$ OF THEM ARE ZERO, AND $n-r-1$ OF THEM ARE ONE; $\underline{q}_1, \underline{q}_2, \dots, \underline{q}_n$ ARE CORRESPONDING NORMED, MUTUALLY ORTHOGONAL EIGENVECTORS, SO WE ALSO KNOW $Q\underline{Q}^T = \underline{Q}^T Q = I$.

NOW, WE CAN GO BACK TO STUDY $\hat{\Sigma}$, WE KNOW

$$\hat{\Sigma} \sim \mathcal{N}_n(0, \sigma^2(I - P))$$

$$\Leftrightarrow \hat{\Sigma} \sim \mathcal{N}_n(0, \sigma^2 Q \Lambda Q^T)$$

$$\Rightarrow Q^T \hat{\Sigma} \sim \mathcal{N}_n(0, \sigma^2 Q^T Q \Lambda Q^T Q)$$

$$\Leftrightarrow Q^T \hat{\Sigma} \sim \mathcal{N}_n(0, \sigma^2 \Lambda) \quad (*)$$

THEN WE OBSERVE THAT

$$\hat{\Sigma}^T \hat{\Sigma} = \hat{\Sigma}^T Q Q^T \hat{\Sigma} = (Q^T \hat{\Sigma})^T \cdot Q^T \hat{\Sigma}$$

WHERE $Q^T \hat{\Sigma}$ IS A COLUMN VECTOR OF RANDOM VARIABLES WITH LENGTH n , AND OBSERVE $(*)$, WE KNOW: $r+1$ MANY COMPONENTS OF $Q^T \hat{\Sigma}$ ARE OF ZERO MEAN, ZERO VARIANCE AND INDEPENDENT TO OTHER COMPONENTS (ZERO COVARIANCE IMPLIES THE UNCORRELATION, GAUSSIAN DISTRIBUTION GIVES THE INDEPENDENCE), THUS WE CAN REGARD THEM AS BEING "DEGENERATED" TO CONSTANT NUMBER 0, THIS GIVES US THE FINAL RESULT:

$$\hat{\Sigma}^T \hat{\Sigma} = (Q^T \hat{\Sigma})^T \cdot Q^T \hat{\Sigma} = \sum_{j=1}^{n-r-1} [Q^T \hat{\Sigma}]_j^2 = \sigma^2 \sum_{j=1}^{n-r-1} \left(\frac{[Q^T \hat{\Sigma}]_j}{\sigma} \right)^2$$

WHERE: $[Q^T \hat{\Sigma}]_j \sim \mathcal{N}_1(0, \sigma^2) \Leftrightarrow \frac{1}{\sigma} [Q^T \hat{\Sigma}]_j \sim \mathcal{N}_1(0, 1)$

$$\Rightarrow \hat{\Sigma}^T \hat{\Sigma} \sim \sigma^2 \chi^2_{(n-r-1)}$$

FROM RESULT 6 WE CAN IMMEDIATELY SEE $\hat{\Sigma}^T \hat{\Sigma}$ CAN BE WRAPPED TO BE A GOOD ESTIMATOR OF σ^2 .

COROLLARY 7 DEFINE $\hat{\sigma}^2 := \frac{\hat{\Sigma}^T \hat{\Sigma}}{n-r-1}$, THEN $\hat{\sigma}^2$ IS AN UNBIASED

COROLLARY 7 DEFINE $s^2 := \frac{\sum \sum}{n-r-1}$, THEN s^2 IS AN UNBIASED ESTIMATOR OF σ^2 .

PROOF SINCE $\sum \sum \sim \sigma^2 \chi^2(n-r-1)$, THEN DIRECTLY

$$E[s^2] = E\left[\frac{\sum \sum}{n-r-1}\right] = \frac{\sigma^2}{n-r-1} \cdot (n-r-1) = \sigma^2 \quad \blacksquare$$

FROM RESULT 6, WE GET $\hat{\beta} \sim N_{r+1}(\beta, \sigma^2 (\mathbf{Z}' \mathbf{Z})^{-1})$. BASE ON THIS WE COULD DO MANY STATISTIC INFERENCES, THE MOST IMPORTANT ONE IS GIVEN AS FOLLOWING:

RESULT 8 LET $\mathbf{Y} = \mathbf{Z}\beta + \varepsilon$, WHERE \mathbf{Z} HAS FULL RANK $r+1$ AND $\varepsilon \sim N_n(0, \sigma^2 I)$. THEN A 100(1- α) PERCENT CONFIDENCE REGION FOR β IS GIVEN BY

IGNORE THE TERMINOLOGY
LIKE "SIMULTANEOUS
CONFIDENCE INTERVALS".

$$(\hat{\beta} - \beta)^T \mathbf{Z} (\hat{\beta} - \beta) \leq (r+1)s^2 F_{r+1, n-r-1}(\alpha)$$

if you don't know it WHERE $F_{r+1, n-r-1}(\alpha)$ IS THE UPPER (100 α)-TH PER CENTILE OF AN F-DISTRIBUTION WITH $r+1$ AND $n-r-1$ DEGREES OF FREEDOM.

ALSO, SIMULTANEOUS 100(1- α) PERCENT CONFIDENCE INTERVALS FOR THE β_i ARE GIVEN BY

$$\hat{\beta}_i \pm \sqrt{\text{Var}(\hat{\beta}_i)} \cdot \sqrt{(r+1) F_{r+1, n-r-1}(\alpha)}, \quad i=0, 1, \dots, r$$

WHERE $\widehat{\text{Var}}(\hat{\beta}_i)$ IS THE DIAGONAL ELEMENT OF ESTIMATED COVARIANCE MATRIX $s^2 (\mathbf{Z}' \mathbf{Z})^{-1}$ CORRESPONDING TO β_i .

PROOF PLEASE REFER TO BOOK "APPLIED MULTIVARIATE STATISTICAL ANALYSIS", WRITTEN BY R. JOHNSON ET. AL. ■

IN FACT, RESULT 8 PLAYS A CRUCIAL ROLE WHEN WE WANT TO KNOW IF THE COEFFICIENT OF PREDICTOR VARIABLE IS STATISTICALLY SIGNIFICANT. OR IN OTHER WORDS, IF THE COEFFICIENT OF A PREDICTOR VARIABLE IS "VERY NEAR" TO ZERO, WE MAY DELETE THIS PREDICTOR VARIABLE FROM THE MODEL, RECONSTRUCT THE DESIGN MATRIX \mathbf{Z} . REDO THE REGRESSION

MAY DELETE THIS PREDICTOR VARIABLE FROM THE MODEL, RECONSTRUCT THE DESIGN MATRIX \bar{Z} , REDO THE REGRESSION.

PRACTITIONERS OFTEN IGNORE THE "SIMULTANEOUS" CONFIDENCE PROPERTY OF INTERVAL ESTIMATES IN RESULT 8. INSTEAD, THEY REPLACE $(r+1) F_{m,n-r-1}(\alpha)$ WITH ONE-AT-A-TIME t VALUE $t_{n-r-1}(\alpha/2)$ AND USE THE INTERVALS

$$\hat{\beta}_i \pm t_{n-r-1}(\alpha/2) \sqrt{\text{Var}(\hat{\beta}_i)}$$

WHEN SEARCHING FOR IMPORTANT PREDICTOR VARIABLES.

MOREOVER, HYPOTHESIS TESTS METHOD COULD BE IMPLEMENTED FOR TESTING HYPOTHESIS LIKE $H_0: \beta_{n_1} = \beta_{n_2} = \dots = \beta_{n_k} = 0$, WHERE $n_1, n_2, \dots, n_k \in \{0, 1, \dots, r\}$. BUT WE DONT INTRODUCE THE DETAIL HERE.

INFERENCES FROM THE ESTIMATED REGRESSION FUNCTION

FROM PRECEDING SECTIONS WE GET OUR ESTIMATED REGRESSION FUNCTION. NOW, SUPPOSE THAT WE GET A NEW DATA

$$\underline{z}_0^f = [1, z_{01}, \dots, z_{0r}]$$

AND WE WANT TO LEARN THE BEHAVIOR OF RESPONSE OF THE ESTIMATED REGRESSION FUNCTION, WE SHALL DENOTE IT AS y_0 .

THEN, HERE ARE TWO SLIGHTLY DIFFERENT STRATEGIES:

1. ESTIMATE THE EXPECTED VALUE OF y_0 .
2. PREDICT y_0 . (MORE UNCERTAIN)

ESTIMATION THE EXPECTED VALUE OF y_0

ACCORDING TO THE REGRESSION MODEL $y = \bar{z}\hat{\beta} + \varepsilon$, IN WHICH $\varepsilon \sim N_n(0, \sigma^2 I)$, THE EXPECTED VALUE OF y_0 IS

$$E[y_0 | \underline{z}_0] = \beta_0 + \beta_1 z_{01} + \dots + \beta_r z_{0r} = \underline{z}_0^f \hat{\beta}$$

ITS LEAST SQUARES ESTIMATION IS $\underline{z}_0^f \hat{\beta}$.

ITS LEAST SQUARE ESTIMATION IS $\underline{\underline{z}}^T \hat{\beta}$.

RESULT 9 FOR THE CLASSICAL LINEAR REGRESSION MODEL WITH $\varepsilon \sim N_n(0, \sigma^2 I)$,
 $\underline{\underline{z}}^T \hat{\beta}$ IS THE UNBIASED LINEAR ESTIMATOR OF $E[y_0 | \underline{\underline{z}}]$ WITH
MINIMUM VARIANCE,

$$\text{Var}(\underline{\underline{z}}^T \hat{\beta}) = \underline{\underline{z}}^T (\underline{\underline{z}} \underline{\underline{z}}^T)^{-1} \underline{\underline{z}} \sigma^2$$

MOREOVER, A $100(1-\alpha)\%$ CONFIDENCE INTERVAL FOR $E[y_0 | \underline{\underline{z}}]$
IS PROVIDED BY

$$\underline{\underline{z}}^T \hat{\beta} \pm t_{n-r-1} \left(\frac{\alpha}{2} \right) \sqrt{(\underline{\underline{z}}^T (\underline{\underline{z}} \underline{\underline{z}}^T)^{-1} \underline{\underline{z}}) \sigma^2}$$

WHERE $t_{n-r-1}(\alpha/2)$ IS THE UPPER $100(\alpha/2)$ -TH PERCENTILE OF
A T-DISTRIBUTION WITH $n-r-1$ DEGREE OF FREEDOM.

FORECASTING A NEW OBSERVATION AT $\underline{\underline{z}}_0$

PREDICTION OF A NEW OBSERVATION, WE CALL IT y_0 , IS MORE UN-
CERTAIN THAN ESTIMATING THE EXPECTED VALUE OF y_0 . ACCORDING
TO THE CLASSICAL REGRESSION MODEL,

$$y_0 = \underline{\underline{z}}_0^T \hat{\beta} + \varepsilon_0$$

WHERE $\varepsilon_0 \sim N_1(0, \sigma^2)$ AND IS INDEPENDENT OF ε , HENCE OF $\hat{\beta}$ AND
 σ^2 . THE ERRORS ε INFLUENCE THE ESTIMATORS $\hat{\beta}$ AND σ^2 THROUGH
THE RESPONSES $\underline{\underline{z}}$, BUT ε_0 DOES NOT.

RESULT 10 GIVEN A CLASSICAL LINEAR REGRESSION MODEL, A NEW
OBSERVATION y_0 HAS THE UNBIASED PREDICTOR

$$\underline{\underline{z}}_0^T \hat{\beta} = \hat{\beta}_0 + \hat{\beta}_1 z_{01} + \dots + \hat{\beta}_r z_{0r}$$

THE VARIANCE OF THE FORECAST ERROR $y_0 - \underline{\underline{z}}_0^T \hat{\beta}$ IS

$$\text{Var}(y_0 - \underline{\underline{z}}_0^T \hat{\beta}) = \sigma^2 (1 + \underline{\underline{z}}_0^T (\underline{\underline{z}} \underline{\underline{z}}^T)^{-1} \underline{\underline{z}}_0)$$

WHEN THE ERRORS ε HAVE A NORMAL DISTRIBUTION, A
100(1- α)% PREDICTION INTERVAL FOR y_0 IS GIVEN BY

$$\underline{\underline{z}}_0^T \hat{\beta} \pm t_{n-r-1}(\alpha/2) \sqrt{\sigma^2 (1 + \underline{\underline{z}}_0^T (\underline{\underline{z}} \underline{\underline{z}}^T)^{-1} \underline{\underline{z}}_0)}$$

$$\hat{y}_0 \hat{\beta} \pm t_{n-r-1}(2) \sqrt{s^2(1 + \hat{z}_0^T (\hat{z}^T \hat{z})^{-1} \hat{z}_0)}$$

WHERE $t_{n-r-1}(2)$ IS THE UPPER 100(2)-TH PERCENTILE OF A t -DISTRIBUTION WITH $n-r-1$ DEGREE OF FREEDOM.

WE CAN OBSERVE THAT: THE PREDICTION INTERVAL FOR \hat{y}_0 IS WIDER THAN THE CONFIDENCE INTERVAL FOR ESTIMATING THE VALUE OF $E[y_0 | \hat{z}_0]$. THE ADDITIONAL UNCERTAINTY IN FORECASTING \hat{y}_0 , WHICH IS REPRESENTED BY EXTRA TERM s^2 IN THE EXPRESSION $s^2(1 + \hat{z}_0^T (\hat{z}^T \hat{z})^{-1} \hat{z}_0)$, COMES FROM THE PRESENCE OF UNKNOWN ERROR ε_0 .

COLLINEARITY

THE COLLINEARITY IS A PHENOMENON IN WHICH ONE PREDICTOR VARIABLE IN A LINEAR REGRESSION MODEL CAN BE UNEASILY PREDICTED FROM THE OTHERS WITH A LARGE DEGREE OF ACCURACY. OR MATHEMATICALLY SPEAKING, THE DESIGN MATRIX IS "NEARLY" NON-FULL-RANK, THEN THE CALCULATION OF $(\hat{z}^T \hat{z})^{-1}$ IS NUMERICALLY UNSTABLE. HERE WE GIVE SOME REMEDIES:

1. **DELETING.** DELETE ONE OF A PAIR OF PREDICTOR VARIABLES THAT ARE STRONGLY CORRELATED.
2. **PRINCIPAL COMPONENT ANALYSIS.** RECALL THAT THIS METHOD PRODUCES NEW UNCORRELATED DATA.
3. **RIDGE REGRESSION.**