

TMA4267 Linear statistical models

Part 2: Linear regression

February 20, 2017

Normal equations

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \text{ where } E(\boldsymbol{\varepsilon}) = \mathbf{0} \text{ and } \text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$$

Which of the following are *the normal equations*?

A $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{H}\mathbf{Y}$

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

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

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

The hat matrix

Design matrix \mathbf{X} has n rows and p linearly independent columns. $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is called the hat-matrix.

Which of the following statements are NOT true?

A $\mathbf{H} = \mathbf{H}^T = \mathbf{H}^2$

B $\text{rank}(\mathbf{H}) = p$

C $\mathbf{H}\mathbf{Y} = \mathbf{Y}$

D $\mathbf{H}(\mathbf{I} - \mathbf{H}) = \mathbf{0}$

Estimator for σ^2

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \text{ where } E(\boldsymbol{\varepsilon}) = \mathbf{0} \text{ and } \text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$$

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

An unbiased estimator for σ^2 is:

A SSE/n

B $\mathbf{Y}^T (\mathbf{I} - \mathbf{H}) \mathbf{Y} / (n - p)$

C $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{Y} / (n - p)$

D $(\mathbf{X}^T \mathbf{X})^{-1} \text{SSE} / n$

Inference about β

$$\mathbf{Y} = \mathbf{X}\beta + \varepsilon \text{ where } \varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$$

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

What are the properties of $\hat{\beta}$?

A Chi-squared distributed with $n - p$ degrees of freedom.

B Chi-squared distributed with p degrees of freedom.

C Multivariate normal with covariance matrix $(\mathbf{I} - \mathbf{H})\sigma^2$.

D Multivariate normal with covariance matrix $(\mathbf{X}^T \mathbf{X})^{-1}\sigma^2$.

$$\text{Happiness} = \text{money} + \text{sex} + \text{love} + \text{work}$$

	Estimate	Std. Error	t value	Pr(> t)
money	0.009578	0.005213	1.837	0.0749
sex	-0.149008	0.418525	-0.356	0.7240
love	1.919279	0.295451	6.496	1.97e-07
work	0.476079	0.199389	2.388	0.0227

Which of the regression coefficient estimates has the largest estimated variance?

A money

B sex

C love

D work

Happiness=money+sex+love+work

The R^2 for the happiness-regression model is 71%. What does that mean?

- A The regression is significant for significance level 71%
- B The regression explains 71% of the variability in the data
- C The estimate for the variance σ^2 is 0.71
- D The covariates have a correlation of 0.71

Happiness

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.072081	0.852543	-0.085	0.9331
money	0.009578	0.005213	1.837	0.0749
sex	-0.149008	0.418525	-0.356	0.7240
love	1.919279	0.295451	6.496	1.97e-07
work	0.476079	0.199389	2.388	0.0227

For which β_j would we reject the null hypothesis $\beta_j = 0$ at significance level 1%?

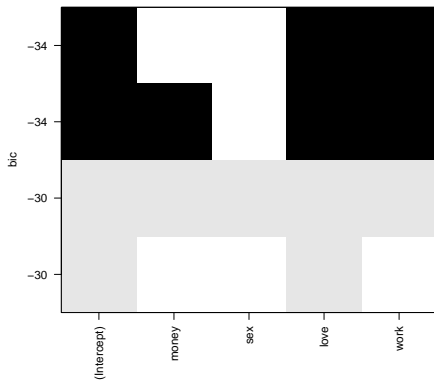
A money

B sex

C love

D work

Best model



Which model does the BIC criterion report to be the best?

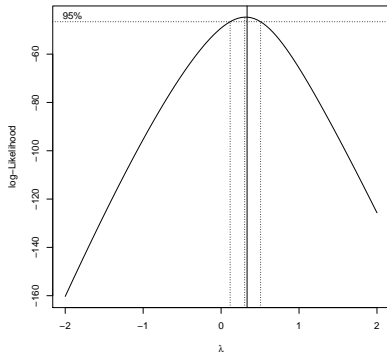
A love+work

B love

C money+love+work

D money+sex+love+work

What is this plot used for?



A Check residuals

B Assess normality of residuals

C Assess linearity

D Find transform of response

Correct?

Are you sure you want to read the correct answers?
Maybe try first? The answers are explained on the
next two slides.

Answers

1. C: The normal equation $(\mathbf{X}^T \mathbf{X}) \hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{Y}$ is before you solve for $\hat{\boldsymbol{\beta}}$.
2. C: The hat matrix is symmetric and idempotent (so A is ok), and has rank p , but the reason for the name of the hat matrix is that it puts the hat on the \mathbf{Y} so $\mathbf{H}\mathbf{Y} = \hat{\mathbf{Y}}$. We know that for symmetric projection matrices the two matrices \mathbf{H} and $(\mathbf{I} - \mathbf{H})$ are orthogonal so the product must be zero.

Answers

3. B: Since SSE has mean $(n - p)\sigma^2$, then $\text{SSE}/(n-p)$ must be an unbiased estimator for σ^2 . We know that $(\mathbf{I} - \mathbf{H})$ projects onto the space orthogonal to the column space of the design matrix, so that must have to do with SSE.
4. D: We know that linear combinations of multivariate normal random vectors are also multivariate normal (so the chi-square is not suitable). The residuals have $(\mathbf{I} - \mathbf{H})$ as part of their covariance matrix, but $\hat{\beta}$ has not.

Answers

5. B: Sex has the largest estimated variance for regression estimate.
6. B: R^2 gives the percent of variability explained.
7. C: only love is significant on level 1%, since this is the only p -value below 0.01 (last column).
8. A: love+work has smallest BIC.
9. D: Box-Cox plot used to find transformation of response.