

STATS FINAL

Study guide

Modeling Basics

$$y = f(x) + \epsilon$$

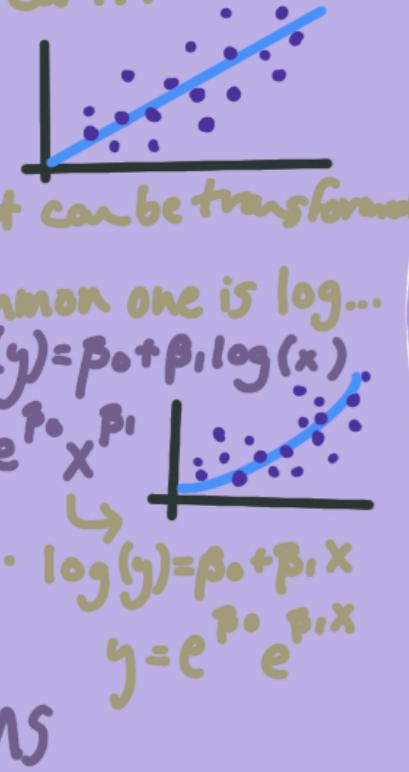
• all models are made up of... response function error

- we can add a variable called a predictor...
- remember... "all models are wrong, but some are useful"

Linear Regression

$$y = \beta_0 + \beta_1 x + \epsilon$$

Model ↓ intercept predictor
↓ response slope



Conditions

Linearity... residual plots

good! even above + below

bad! shaped and fan

Constant Variance...

good! even above + below

bad! not even fan shape

Normality...

good! close to line straight

bad! not on line shaped

Independence

Knowing one error won't help us find another

Mean

errors centered at mean

help us find another

guaranteed by LSR

Outliers

not exactly a condition, but still important to check.

Outliers are points far from the regression line

Influential points are points that pull the regression fit.

Studentized Res. test for outliers > 2.5

susy bad!

Leverage tests for influential points > $\frac{4}{n}$

susy bad!

Cook's Distance combines both.

>.5 > 1.0

susy bad!

Interpretations

$\beta_0 \rightarrow$ intercept.

"when the predictors in our model predicts a response of β_0 "

$\beta_1 \rightarrow$ slope

"for every 1 change in our predictor our model predicts a β_1 change in response"

β_0 (logtrans) \rightarrow

"when the predictor is 2 our expected response is β_0 "

β_1 (logtrans) \rightarrow

"changing our predictor by e^{β_1} our model predicts a change in response of e^{β_1} "

Uses

Hypothesis Test

is it possible there is no relationship between predictor and response

$H_0: \beta_1 = 0$ $H_1: \beta_1 \neq 0$

p-value, if small reject H_0

Confidence Interval

95% conf. interval for true slope...

"we are 95% sure the true slope is between these values"

Conf. interval for true mean response for predictor x_k

"we are 95% con. that the true mean resp. for x_k is ..."

What would I expect the next response for x_k to be...

"we are 95% con. that the next response for a predictor x_k will be between these values."

Variability Explained by Model

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

bad \rightarrow good

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$$

Model ↓ intercept predictor

linear...

but can be transformed...

common one is log...

$\log(y) = \beta_0 + \beta_1 \log(x)$

$y = e^{\beta_0 + \beta_1 x}$

or... $\log(y) = \beta_0 + \beta_1 x$

$y = e^{\beta_0 + \beta_1 x}$

Conditions

Linearity... residual plots

good! even above + below

bad! shaped and fan

Constant Variance...

good! even above + below

bad! not even fan shape

Normality...

good! close to line straight

bad! not on line shaped

Independence

Knowing one error won't help us find another

Mean

errors centered at mean

help us find another

Guaranteed by LSR

Outliers

not exactly a condition, but still important to check.

Outliers are points far from the regression line

Influential points are points that pull the regression fit.

Studentized Res. test for outliers > 2.5

susy bad!

Leverage tests for influential points > $\frac{4}{n}$

susy bad!

Cook's Distance combines both.

>.5 > 1.0

susy bad!

Interpretations

$\beta_0 \rightarrow$ intercept.

"when the predictors in our model predicts a response of β_0 "

$\beta_1 \rightarrow$ slope

"for every 1 change in our predictor our model predicts a β_1 change in response"

β_0 (logtrans) \rightarrow

"when the predictor is 2 our expected response is β_0 "

β_1 (logtrans) \rightarrow

"changing our predictor by e^{β_1} our model predicts a change in response of e^{β_1} "

Uses

Hypothesis Testing

perform indiv. t-test to see if single coef is useful

same as intep. as single

perform nested F test on full category...

$H_0: \beta_2 = \beta_3 = 0$

$H_1: \beta_2 \text{ or } \beta_3 \neq 0$

$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x_{cat}$

vs

$y = \beta_0 + \beta_1 x$

"we have statistically sig. evidence that the true β_2/β_3 is not 0 and the categorical variable has a relation w/ the response!"

Confidence Intervals

same as last time

Model Selection

remember \rightarrow short and sweet.

Balance high variability explained (high R^2) w/ a small-ish model.

Can also use AIC

automated model selection methods are no substitute for thinking

Y-axis

Good Luck

Y-axis

Good Luck