

STAT 213

Model Selection

Colin Reimer Dawson

Oberlin College

March 26, 2018

1 / 25

Notes

Outline

Review: Simpson's Paradox

Added Variable Plots

Model Selection
 Penalized Fit Measures

2 / 25

Notes

Outline

Review: Simpson's Paradox

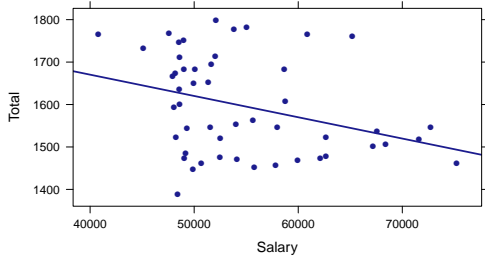
Added Variable Plots

Model Selection
 Penalized Fit Measures

3 / 25

Notes

SATs and Teacher Salary

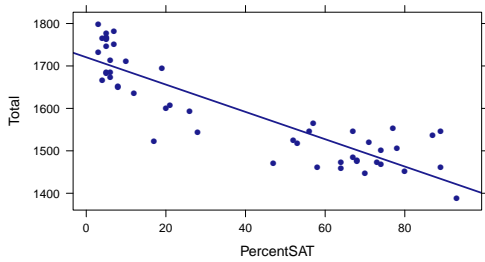


Overall, states that pay teachers more have lower mean SAT scores.

4 / 25

Notes

SATs and Teacher Salary

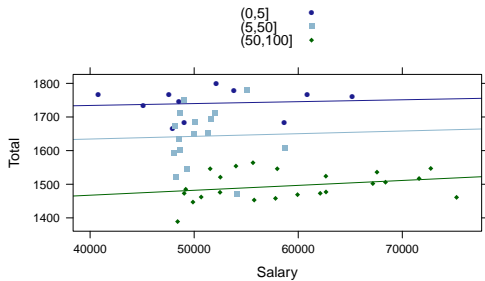


Overall, states with higher participation rates have lower mean SAT scores.

5 / 25

Notes

Controlling for Participation



Controlling for participation, the more teachers are paid, the higher the mean score.

6 / 25

Notes

Controlling for Participation

```
m.salary.percent <- lm(Total ~ Salary + PercentSAT, data = SATdata)
m.salary.percent %>% summary() %>% coefficients() %>% round(3)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1589.007	58.471	27.176	0.000
Salary	0.003	0.001	2.295	0.026
PercentSAT	-3.553	0.278	-12.756	0.000

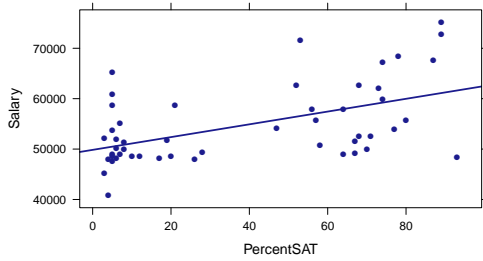
In the full model the Salary coefficient has the opposite sign from the simple model! This is an instance of **Simpson's Paradox**

7 / 25

Notes

What's Going On?

```
xyplot(Salary ~ PercentSAT, data = SATdata, type = c("p", "x"))
```



PercentSAT was a **confounding variable**

8 / 25

Notes

Outline

Review: Simpson's Paradox

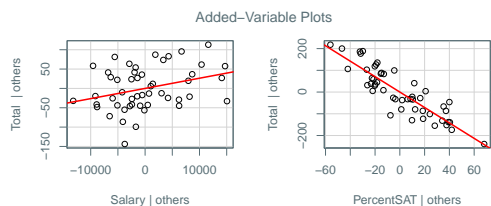
Added Variable Plots

Model Selection
Penalized Fit Measures

9 / 25

Notes

Visualizing "Added Value"



To see the unique contribution of a predictor, X_k , we can plot two sets of residuals against each other

1. Residuals from model predicting Y from other predictors
2. Residuals from model predicting X_k from other predictors

Notes

Added Variable Plots

- Also called "partial regression plots"
- Idea: Capture the *nonredundant* information provided by X_k
- How much information does X_k provide about Y that we couldn't have already predicted?

Notes

Example: Quiz and Exam Scores

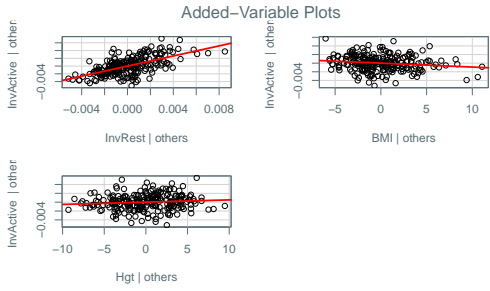
```
m.both <- lm(Final ~ Quiz + Midterm, data = Scores)
avPlots(m.both)
```



Notes

Example: Pulse Rates

```
m.pulse <- lm(InvActive ~ InvRest + BMI + Hgt, data = Pulse2)
avPlots(m.pulse)
```

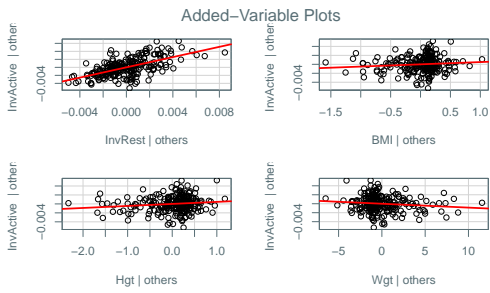


13 / 25

Notes

Example: Pulse Rates (With Multicollinearity)

```
m.pulse2 <- lm(InvActive ~ InvRest + BMI + Hgt + Wgt, data = Pulse2)
avPlots(m.pulse2)
```



Note the differences in x axis range!

14 / 25

Notes

Outline

Review: Simpson's Paradox

Added Variable Plots

Model Selection
Penalized Fit Measures

15 / 25

Notes

So many models...

- How to decide among all these models?
 1. Understand the subject area! Build sensible models.
 2. Nested F -tests to make *targeted* comparisons of interest
 3. Model quality measures

16 / 25

Notes

ASSESS: Coefficient of Determination

Can we use R^2 to identify the best model?

$$\text{As before, } R^2 = \frac{SS_{Model}}{SS_{Total}} = 1 - \frac{SS_{Error}}{SS_{Total}}$$

17 / 25

Notes

What Happens if We Add Useless Predictors?

Simulation Handout

18 / 25

Notes

What Makes a "Good" Model?

Fit

- High R^2
- Small SSE (equivalent)

Validity

- Satisfies stated conditions
- Strong evidence for predictors
- Simple (Parsimonious)
- Generalizes to new data

19 / 25

Notes

Why Does Parsimony Matter?

Don't we just care about good predictions?

Not exclusively...

- We also use models to *understand* the world (harder with more complexity)

And even so...

- We really care about making predictions for data we *haven't seen yet*.

20 / 25

Notes

Balancing Fit and Parsimony

- R^2 can only go up as we add a predictor, because at worst, we can choose $\beta_{K+1} = 0$ and get the same SSE. Usually we can pick coefficients to do somewhat better.
- Would like to "penalize" unnecessary predictors.

22 / 25

Notes

Adjusted R^2 (Three equivalent formulas)

$$\begin{aligned}
 R_{adj}^2 &= 1 - \frac{SS_{Error}/(N - K - 1)}{SS_{Total}/(N - 1)} \\
 &= 1 - \frac{\hat{\sigma}_\varepsilon^2}{s_Y^2} \\
 1 - R_{adj}^2 &= (1 - R^2) \frac{df_{Total}}{df_{Error}}
 \end{aligned}$$

23 / 25

Notes

Criteria to "score" models

1. high R^2 /low SSE/low $\hat{\sigma}_\varepsilon^2$: always favors more complex models
2. Adj. R^2 : balances fit and complexity
3. Mallows's C_p / Akaike Information Criterion (AIC): estimates mean squared prediction error based on $\hat{\sigma}_\varepsilon^2$ from a "full" model
4. Out-of-sample predictive accuracy

24 / 25

Notes

Mallows's C_p / AIC

Two measures that reduce to the same thing in the case of MLR with independent, equal variance, Normal residuals. For a "reduced" model with $p_{reduced}$ total parameters (including the intercept) which is nested in a "full" model with p_{full} parameters, both fit using n observations:

$$C_p = \frac{SSE_{reduced}}{MSE_{full}} + 2p_{reduced} - n \quad (1)$$

$$= p_{reduced} + \frac{SS_{nested}}{MSE_{full}} \quad (2)$$

where SS_{nested} is the (positive) difference between the sums of squares of the two models (as for a nested F test)

Should we prefer larger or smaller values?

25 / 25

Notes
