Understanding regression models

ST552 Lecture 10

Charlotte Wickham 2019-01-30

- Lecture: mathematical strategies for understanding models
- Lab: understanding models through visualization

We've talked about the machinery to perform:

- t-tests, t-based confidence intervals for individual βs and linear combinations of βs
- F-tests for hypotheses about many βs

But what can we do with this machinery?

- 1. Understand a model in the context of a problem
- 2. Define a set of models to answer questions of interest

Focus on #1 to gain intuition in how to approach #2.

Example

LA Dodgers (baseball team) sometimes give out "bobbleheads" at home games. They are curious if this increases attendance at games.

81 games in 2012 season, 11 of which bobbleheads were given out. Have measurements on:

- attendance at game (number of people)
- day of the week the game was played Mon, Twe, Wed ...
- some other variables, that are probably important, but we will ignore for now

(Inspired by Chapter 2 in Modeling Techniques in Predictive Analytics: Business Problems and Solutions with R. Get data from http:

//www.informit.com/promotions/modeling-techniques-in-predictive-analytics-141183 if you interested.)

attendance_i =
$$\beta_0 + \beta_1 1$$
{bobblehead YES}_i+
 $\beta_2 1$ {Tue}_i + $\beta_3 1$ {Wed}_i + $\beta_4 1$ {Thu}_i+
 $\beta_5 1$ {Fri}_i + $\beta_6 1$ {Sat}_i + $\beta_7 1$ {Sun}_i + ϵ_i

What does this model say about the relationship between attendance and whether bobbleheads are given out and day of the week? Whese is Monday?

6

home games

variable: a measurement made on the observational units.

E.g. bobblehead (yes/no) and day of week (mon/tue/wed/thu/fri/sat/sun).

term: a column of the design matrix

E.g. bobblehead, $1{Fri}$

do we count htercept? Yos

P=# parameters = # terms

Ask about the effect of a variable?

What does the model say about the mean response when a variable is varied, holding all other variables constant.

Categorical variable: Find the mean response for each level and compare.

Continuous variable: Find the change in mean response if the variable increases by 1 unit.

The variable has two levels yes no

We'll find E (attendence bobblehead = Yes) and E (attendence bobblehead = No) then compare them.

$$E(\text{attendence}|\text{ bobblehead} = Yes) \qquad \text{beld constant}$$

$$= \mathcal{B}_{0} + \mathcal{B}_{1} \mathbf{1} + \mathcal{B}_{2} \mathbf{t} + \mathcal{B}_{5} \mathbf{\omega} + \mathcal{B}_{4} \mathbf{h} + \mathcal{B}_{5} \mathbf{f} + \mathcal{B}_{6} \mathbf{s} + \mathcal{B}_{7} \mathbf{c}$$

$$E(\text{attendence}|\text{ bobblehead} = \text{No})$$

$$= \mathcal{B}_{0} + \mathcal{B}_{2} \mathbf{t} + \mathcal{B}_{5} \mathbf{\omega} + \mathcal{B}_{4} \mathbf{h} + \mathcal{B}_{5} \mathbf{f} + \mathcal{B}_{6} \mathbf{s} + \mathcal{B}_{7} \mathbf{u}$$

E (attendence | bobblehead = Yes) - E (attendence | bobblehead = No) $= \beta_{i}$

- For a fixed day of the week, the model predicts the mean attendance increases by β₁ when bobbleheads are given out.
- The model predicts the mean attendance when bobbleheads are given is β_{1×} more than when bobbleheads aren't given, after accounting for day of the week.

If β_1 is zero, then bobbleheads don't have an effect on attendance.

We could answer the questions:

- Is the mean attendance higher when bobbleheads are given out? t-test on β₁ = 0.
- How much higher is the mean attendance higher when bobbleheads are given out? Confidence interval on β₁

But this is an observational study, so we need to be careful with our language!

OK "It is estimated distributing bobbleheads is **associated** with an increased mean attendance of XX". "Special" given at "special" games Not OK "It is estimated distributing bobbleheads increases the mean attendance by XX".

What does model say about the expected attendance on Monday? What does model say about the expected attendance on Tuesday? What does model say about the expected attendance on Wednesday? Your turn: Day of the week

E (attend Day of week = Mon) = Bot B, 1 { bobblehead } = Bo+B, b+B21{Tre}; + O $M_{\delta n} := \beta_0 + \beta_1 b$ pz: difference in mean attend Tue: Wed: Bo + Bib + Bz Bo + B(b + B3 between Mon & Tue holding everything 14

If $\beta_2 = \beta_3 = \beta_4 = \beta_5 = \beta_6 = \beta_7$ is zero, then day of the week doesn't have an effect on attendance.

What tools could we use to answer the questions:

- Does day of the week have an effect on mean attendance (after accounting for the bobblehead promotion)? F-test
- How much does the mean attendance differ between Friday and Saturday? testing: competing models F-lest t-test flo: BS-B6=0
 CF: BS-B6

In lab today:

weight_i =
$$\beta_0 + \beta_1 1 \{ \text{male} \}_i + \beta_2 \text{height}_i + \epsilon_i$$

What is the effect of height?

Compare E (weight| height = h) to E (weight| height = h + 1), holding other variables constant.

Interactions: describe situations where the effect of one variable depends on the level of another explanatory variable.

E.g.

weight_i = $\beta_0 + \beta_1 1 \{ \text{male} \}_i + \beta_2 \text{height}_i + \beta_3 (1 \{ \text{male} \} \times \text{height})_i + \epsilon_i$ The same strategy will work. Sex hgt in × looks like male 170 (170) female 160 (0) Semale 160 (165) 17

wyti = Bo+B, 12 male], + Bz hgt, 1 β3 { 1?male3×hgt3; + ε. What is the effect of height ? E (wat | hat=h, 12male3=m3 = Bo+B, m+Bzh+Bzmh $E(wat|hat=h+l,m) = \beta_0 + \beta_1 m + \beta_2(h+l) + \beta_3 m(h+l)$ différence = = Bot Bim + Bzh + Bz + Bzmh + Bzm female = Bz male = Bz + Bz $effect = (B_2 + B_3 m)$ height