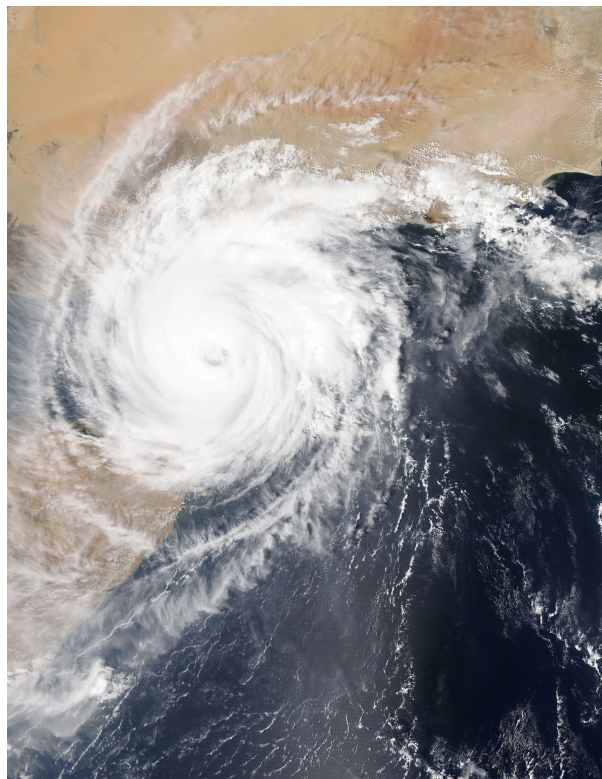


Lecture 2.3 - Developing regression wisdom - sample solution

Student

2025-11-07



Interpretation steps

Interpretation goals

Interpreting the output of regression tables is an essential skill for a statistical analyst. Being able to tell if your regression coefficients indicate a meaningful relationship as well as reading

the diagnostic information found in a regression table is key to model building and statistical literacy. We will develop these skills in three exercises.

- Learning objectives:
 - Practice interpreting slopes
 - Practice interpreting intercepts
 - Develop skills to understand statistical magnitude
 - Understand under what domains models are valid

Interpretation process

These steps are steps that every statistical modeler should take when analyzing data.

1. Before starting, write a sentence or two on your expectation of the relationship – what direction will it be, how strongly will the two be related, and does the predictor variable have a large or small impact on the response variable.

Let's take `LF.WindsKPH` as the predictor and `BaseDam2014` as the response variable

2. What are the units of the predictor and response variable?

Speed in KPH and US 2014 dollars.

3. Draw a scatterplot of the relationship between the two variables. Does it appear that one variable needs to be reexpressed? If so, find the best reexpression possible (the data may not look very pretty even after reinterpretation, your goal is to make the relationship as linear as possible)

Yes, does look like it could benefit from a reexpression.

Log transformed damage looks (somewhat) better

4. Conduct and view the results of your regression. Store it in a nice table using the `modelsummary()` command.
5. Interpret the slope of the regression line – does a one unit change in the predictor variable lead to a large or small change in the response variable?

A one-unit change in windspeed (1 km/h windspeed increase) predicts a 0.03 log change in damage. In casual terms, we can say that it results in about a 3% increase in predicted damage.

6. What are reasonable values for the predictor variables? For what range of x would our best fit line have meaning and when would it not?

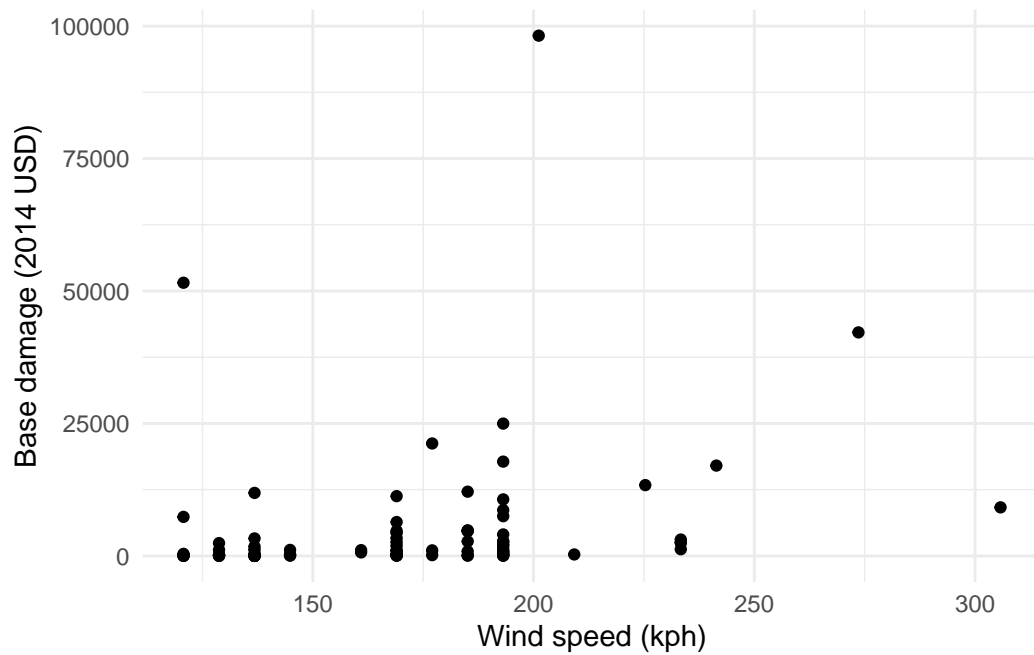


Figure 1: Wind speed vs. damage

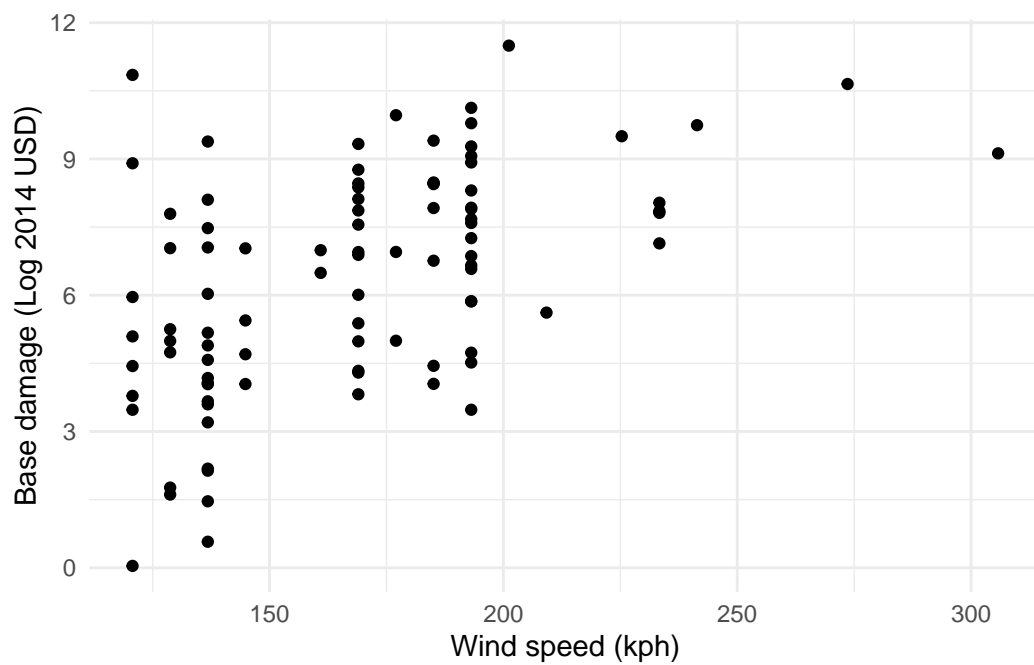


Figure 2: Wind speed vs. log damage

	(1)
(Intercept)	0.419 (1.038)
LF.WindsKPH	0.035 (0.006)
Num.Obs.	94
R2	0.272
R2 Adj.	0.264
AIC	410.4
BIC	418.0
Log.Lik.	−202.191
RMSE	2.08

Reasonable values (in this context) would be something like 80 km/hr to 200 km/hr. Lower speed probably won't produce much (any?) damage and higher is not physically possible.

7. Interpret the intercept. Does it have meaning in this case?

It is the amount of damage when the windspeed is zero, which makes no sense.

8. Solve for the predicted value of the response variable at Q1 and Q3 of the predictor variable. Does moving from Q1 to Q3 of the predictor variable result in a large change (in your opinion) in the response variable?

```
intercept <- mod$coefficients[1]
slope <- mod$coefficients[2]
predq1 <- quantile(hurricanes$LF.WindsKPH, 0.25)
predq3 <- quantile(hurricanes$LF.WindsKPH, 0.75)

print(unname(intercept + slope*predq1))
```

```
[1] 5.250019
```

```
print(unname(intercept + slope*predq3))
```

```
[1] 7.239344
```

A two log-unit change in predicted damage is relatively large given the Q1 to Q3 of the response variable goes from roughly 4 to 8.

9. Draw a box plot by hand of the residuals based on the regression results you can see from `summary(<your model name>)`. Does the boxplot suggest any problems? What are the units?

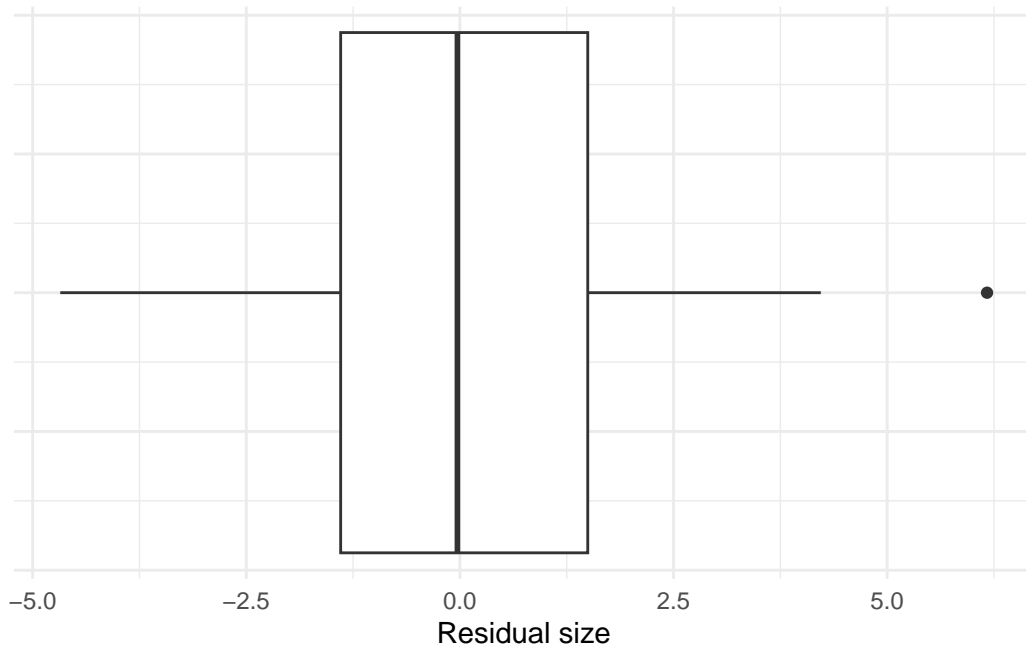


Figure 3

While not drawn by hand, we can see the residuals look roughly symmetrical with only one significant outlier.

10. Interpret the residual standard error from these same results – how big or small is it?

```
summary(mod)$sigma
```

```
[1] 2.101771
```

```
exp(summary(mod)$sigma)
```

```
[1] 8.180643
```

The standard error indicates the average miss is about 2.1 log units (or about 8.1 unlogged units). Given that the reasonable range of logged damage goes from about 4 to 8, that indicates a rather large ‘average’ miss.

11. Make a scatterplot of the fitted values (\hat{y}) vs. the residuals with a line at the 0 value of the y axis. Interpret this plot.

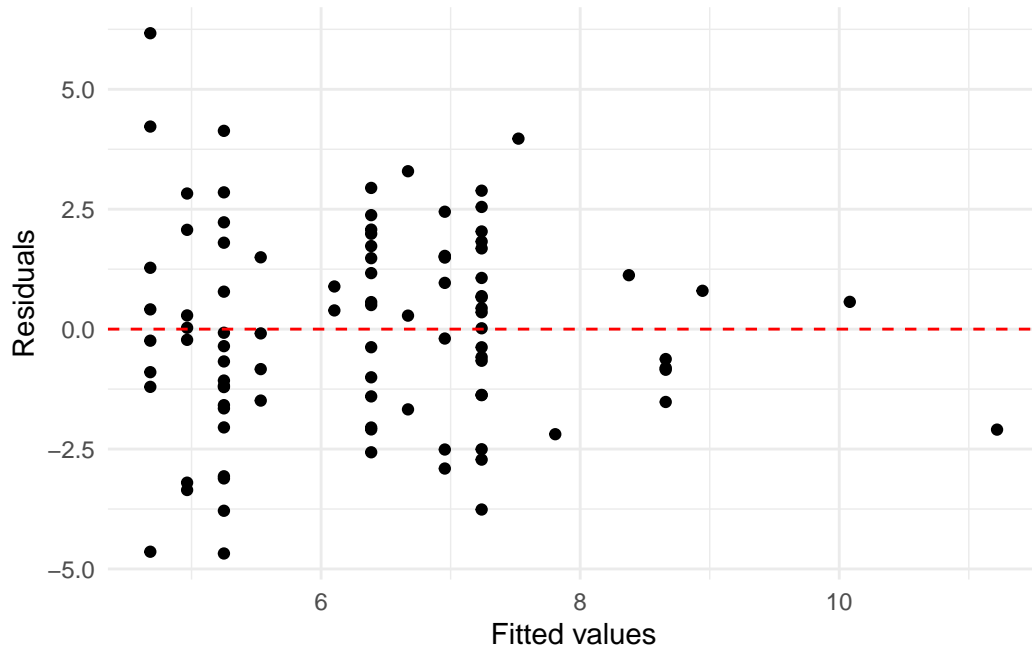


Figure 4: Residuals plot of the model: windspeed vs. log 2014 USD damage

Overall, there does not appear to be any non-linearities in the residuals. However, there does appear to be a significant Plot Thickens condition problem.

11. Interpret the R^2 Specifically, how much of the variance in the response variable is this model explaining?

```
summary(mod)$r.squared
```

```
[1] 0.2719493
```

This moderately low level of R squared is not too surprising - storms cause more damage the higher the wind speed but there are a lot of other factors as well, including whether it hit a populated area or not.

13. What is some more information you'd like to know before making a final conclusion about the relationship between the two variables? Do you think there are any lurking variables?

Probably the main lurking variable is the population size of the area the storm hit. And also the diameter of the storm - a small diameter storm has less possibility of causing damage.

Real data

Hurricane Data

Variable definitions:

- **Name** - Hurricane name
- **Year** - Numeric
- **LF.WindsMPH** - Maximum sustained windspeed (≥ 1 minute) to occur along the US coast. Prior to 1980, this is estimated from the maximum windspeed associated with the Saffir-Simpson index at landfall. If 2 or more landfalls, the maximum is taken
- **LF.PressureMB** - Atmospheric pressure at landfall in millibars. If 2 or more landfalls, the minimum is taken
- **LF.times** - Number of times the hurricane made landfall
- **BaseDamage** - Property damage (in millions of dollars for that year)
- **NDAM2014** - Damage, had hurricane appeared in 2014
- **AffectedStates** - Affected states (2-digit abbreviations), pasted together
- **firstLF** - Date of first landfall
- **deaths** - Number of continental US direct and indirect deaths
- **mf** - Gender of name
- **LF.WindsKPH** - Maximum sustained windspeed expressed in kilometers/hr

Investigate

Using the 13 steps described above, investigate the relationship between two variables that you think should be related. Write up your response to each of the fourteen points listed above. If you finish early, pick another predictor variable that you think may be related to the response variable and repeat the process.