Practical challenge 3

Model crafting Foundations of Data Science—DJW—2018/2019

These questions are optional. They are not intended for supervision, and they are not examinable material. But it is HIGHLY RECOMMENDED that you attempt question 1.

If you want guidance, you can (i) ask your fellow students for help in the Moodle forum, (ii) come to the practical classes in the second half of term and ask demonstrators, (iii) wait until model solutions are released, after the practical classes.

Question 1. Find a 95% confidence interval for the rate of temperature increase at the Cambridge weather station, using the model

$$\mathsf{temp} \approx \alpha + \beta_1 \sin(2\pi \mathsf{t}) + \beta_2 \cos(2\pi \mathsf{t}) + \gamma \mathsf{t}. \tag{1}$$

- url = 'https://teachingfiles.blob.core.windows.net/founds/climate.csv'
- 2 climate = pandas.read_csv(url)
- 3 df = climate.loc[climate.station='Cambridge']
- 4 t = df.yyyy + (df.mm-1)/12
- 5 temp = (df.tmin + df.tmax)/2

Question 2. In this question, you will investigate three refinements of model (1). Answer each part separately; you don't need to combine your findings between the parts. This is a question about designing features for linear models; you don't need to find confidence intervals or conduct hypothesis tests, or devise fancy non-linear probabilistic models.

(a) The Met Office says that it's not just a question of temperature increase, but that the temperatures are getting more extreme. Fit a model to Cambridge temperatures, in which the summer high and the winter low are both changing linearly, but at different rates. What are these rates?

https://www.bbc.co.uk/news/science-environment-46064266

- (b) Perhaps the annual cycle isn't a perfect sin curve. Fit a model to Cambridge temperatures, in which the annual cycle is an arbitrary shape, using one-hot coding as in page 66 of lecture notes. Plot a bar chart of the annual cycle, one bar per month, and superimpose the fitted sin curve from model (1).
- (c) For every station separately, using temperature data from 1960 onwards, fit the model (1). Plot a scatter plot of the α coefficient as a function of latitude, with one dot per station. Plot three more scatter plots, to show how α and γ coefficients depend on latitude and altitude. Propose and fit a single model for the entire dataset (from 1960 onwards), incorporating your findings about α and γ .