

Statistics of Big Data, Spring 2015  
Warmup homework exercise

Due date: 22 March 2015

This exercise relies on the 2009 Nature paper<sup>1</sup> describing Google Flu Trends and the 2014 Science paper<sup>2</sup> discussing the major errors this model made starting in 2011. Both papers are available from the class home page.

1. Read the Nature paper, identify at least two flaws in the statistical methodology, and explain how each of these aspects could have been better addressed. Make sure your answers are specific, accurate and concise (and as mathematical as possible), rather than vague statements. The goal here is not to propose completely different approaches and ideas, but to concentrate on how they could have implemented their chosen approach better. Some of the aspects you may want to consider for this are:
  - (a) The internal four-fold cross validation in fitting each one of the  $5 \cdot 10^6 \cdot 9$  models : is it necessary? Could the same goal have been accomplished with a simpler approach?
  - (b) The methodology for combining the 36  $Z$  scores into one score: is there an obvious better / more powerful approach? Possibly combined with the solution of the previous item.
  - (c) Given that they want to predict the actual percentage in the future, would the use of other evaluation measures in the final step (selecting the number of terms) be more appropriate? Which and why?
2. Beyond general disparaging comments, the Science paper seems to attribute the failure of GFT to two main flaws:
  - (a) Failure to take advantage of CDC data in modeling and failure to compare to baseline based on CDC data auto-regression.
  - (b) “Algorithm dynamics” which seems to relate to how Google is collecting the data and how users are using the search engine. Because these are changing over time, the model no longer predicts well and is consistently overestimating.

Explain each point briefly in concise and mathematical terms. Especially for the second point (“dynamics”) explain which more general phenomenon is being described here that applies to all data that are collected and analyzed over time.

3. Of the points in question 1, and the two points in question 2, which one do you think is the most responsible for the extreme deterioration of GFT performance over time? Explain clearly and briefly. As much as possible use concrete statistical terms like “bias”, “variance”, “stationarity” in their correct meaning in explaining your answer.

---

<sup>1</sup><http://www.nature.com/nature/journal/v457/n7232/full/nature07634.html>

<sup>2</sup><http://www.sciencemag.org/content/343/6176/1203>