# LOCKDOWN SCEPTICS
## STAY SCEPTICAL. CONTROL THE HYSTERIA. SAVE LIVES.

# Code Review of Ferguson's Model

*6 May 2020. Updated 10 May 2020.*

by Sue Denim (not the author's real name)

[Please note: a follow-up analysis is now available here.]

Imperial finally released a derivative of Ferguson's code. I figured I'd do a review of it and send you some of the things I noticed. I don't know your background so apologies if some of this is pitched at the wrong level.

**My background.** I have been writing software for 30 years. I worked at Google between 2006 and 2014, where I was a senior software engineer working on Maps, Gmail and account security. I spent the last five years at a US/UK firm where I designed the company's database product, amongst other jobs and projects. I was also an independent consultant for a couple of years. Obviously I'm giving only my own professional opinion and not speaking for my current employer.

**The code.** It isn't the code Ferguson ran to produce his famous Report 9. What's been released on GitHub is a heavily modified derivative of it, after having been upgraded for over a month by a team from Microsoft and others. This revised codebase is split into multiple files for legibility and written in C++, whereas the original program was "a single 15,000 line file that had been worked on for a decade" (this is considered extremely poor practice). A request for the original code was made 8 days ago but ignored, and it will probably take some kind of legal compulsion to make them release it. Clearly, Imperial are too embarrassed by the state of it ever to release it of their own free will, which is unacceptable given that it was paid for by the taxpayer and belongs to them.

**The model.** What it's doing is best described as "SimCity without the graphics". It attempts to simulate households, schools, offices, people and their movements, etc. I won't go further into the underlying assumptions, since that's well explored elsewhere.

**641**

**Non-deterministic outputs.** Due to bugs, the code can produce very different results given identical inputs. They routinely act as if this is unimportant.

This problem makes the code unusable for scientific purposes, given that a key part of the scientific method is the ability to replicate results. Without replication, the findings might not be real at all – as the field of

psychology has been finding out to its cost. Even if their original code was released, it's apparent that the same numbers as in Report 9 might not come out of it.

Non-deterministic outputs may take some explanation, as it's not something anyone previously floated as a possibility.

The documentation says:

> *The model is stochastic. Multiple runs with different seeds should be undertaken to see average behaviour.*

"Stochastic" is just a scientific-sounding word for "random". That's not a problem if the randomness is intentional pseudo-randomness, i.e. the randomness is derived from a starting "seed" which is iterated to produce the random numbers. Such randomness is often used in Monte Carlo techniques. It's safe because the seed can be recorded and the same (pseudo-)random numbers produced from it in future. Any kid who's played Minecraft is familiar with pseudo-randomness because Minecraft gives you the seeds it uses to generate the random worlds, so by sharing seeds you can share worlds.
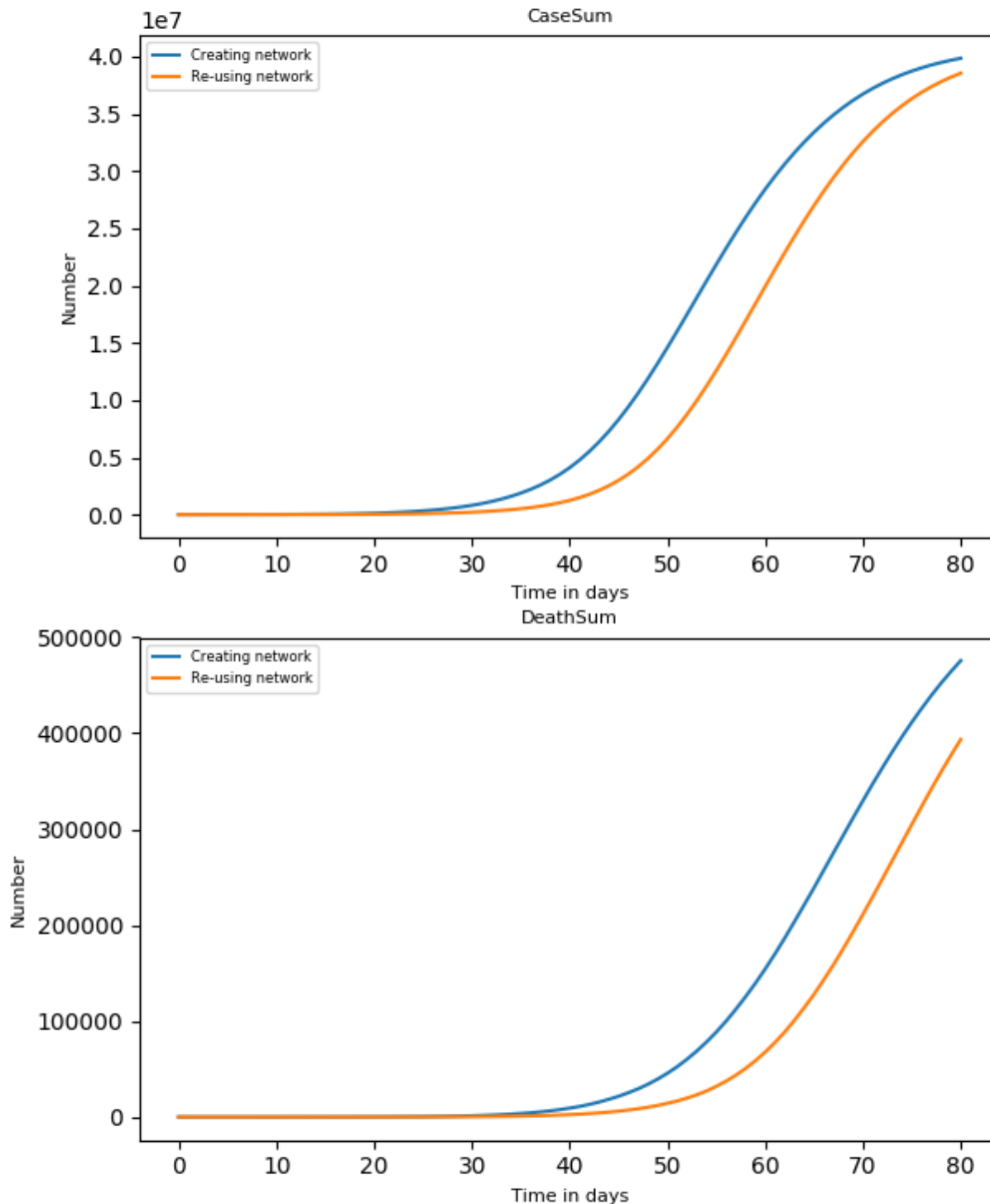
Clearly, the documentation wants us to think that, given a starting seed, the model will always produce the same results.

Investigation reveals the truth: the code produces critically different results, even for identical starting seeds and parameters.

I'll illustrate with a few bugs. In issue 116 a UK "red team" at Edinburgh University reports that they tried to use a mode that stores data tables in a more efficient format for faster loading, and discovered – to their surprise – that the resulting predictions varied by around 80,000 deaths after 80 days:

**641**

**CaseSum**

Legend: Creating network (blue), Re-using network (orange). Y-axis: Number (×1e7), X-axis: Time in days.

**DeathSum**

Legend: Creating network (blue), Re-using network (orange). Y-axis: Number, X-axis: Time in days.

**641**

That mode doesn't change anything about the world being simulated, so this was obviously a bug.

The Imperial team's response is that it doesn't matter: they are "aware of some small non-determinisms", but "this has historically been considered acceptable because of the general stochastic nature of the model".

Note the phrasing here: Imperial know their code has such bugs, but act as if it's some inherent randomness of the universe, rather than a result of amateur coding. Apparently, in epidemiology, a difference of 80,000 deaths is "a small non-determinism".

Imperial advised Edinburgh that the problem goes away if you run the model in single-threaded mode, like they do. This means they suggest using only a single CPU core rather than the many cores that any video game would successfully use. For a simulation of a country, using only a single CPU core is obviously a dire problem – as far from supercomputing as you can get. Nonetheless, that's how Imperial use the code: they know it breaks when they try to run it faster. It's clear from reading the code that in 2014 Imperial tried to make the code use multiple CPUs to speed it up, but never made it work reliably. This sort of programming is known to be difficult and usually requires senior, experienced engineers to get good results. Results that randomly change from run to run are a common consequence of thread-safety bugs. More colloquially, these are known as "Heisenbugs".

But Edinburgh came back and reported that – even in single-threaded mode – they still see the problem. So Imperial's understanding of the issue is wrong.  Finally, Imperial admit there's a bug by referencing a code change they've made that fixes it. The explanation given is "It looks like historically the second pair of seeds had been used at this point, to make the runs identical regardless of how the network was made, but that this had been changed when seed-resetting was implemented". In other words, in the process of changing the model they made it non-replicable and never noticed.

Why didn't they notice? Because their code is so deeply riddled with similar bugs and they struggled so much to fix them that they got into the habit of simply averaging the results of multiple runs to cover it up… and eventually this behaviour became normalised within the team.

In issue #30, someone reports that the model produces different outputs depending on what kind of computer it's run on (regardless of the number of CPUs). Again, the explanation is that although this new problem "will just add to the issues" …  "This isn't a problem running the model in full as it is stochastic anyway".

Although the academic on those threads isn't Neil Ferguson, he is well aware that the code is filled with bugs that create random results. In change #107 he authored he comments: "It includes fixes to InitModel to ensure deterministic runs with holidays enabled".  In change #158 he describes the change only as "A lot of small changes, some critical to determinacy".

641

Imperial are trying to have their cake and eat it.  Reports of random results are dismissed with responses like "that's not a problem, just run it a lot of times and take the average", but at the same time, they're fixing such bugs when they find them. They know their code can't withstand scrutiny, so they hid it until professionals had a chance to fix it, but the damage from over a decade of amateur hobby programming is so extensive that even Microsoft were unable to make it run right.

**No tests.** In the discussion of the fix for the first bug, Imperial state the code used to be deterministic in that place but they broke it without noticing when changing the code.

Regressions like that are common when working on a complex piece of software, which is why industrial software-engineering teams write automated regression tests. These are programs that run the program with varying inputs and then check the outputs are what's expected. Every proposed change is run against every test and if any tests fail, the change may not be made.

The Imperial code doesn't seem to have working regression tests. They tried, but the extent of the random behaviour in their code left them defeated. On 4th April they said: "However, **we haven't had the time** to work out a scalable and maintainable way of running the regression test in a way that allows a small amount of variation, but doesn't let the figures drift over time."

Beyond the apparently unsalvageable nature of this specific codebase, testing model predictions faces a fundamental problem, in that the authors don't know what the "correct" answer is until long after the fact, and by then the code has changed again anyway, thus changing the set of bugs in it. So it's unclear what regression tests really mean for models like this – even if they had some that worked.

**Undocumented equations.** Much of the code consists of formulas for which no purpose is given. John Carmack (a legendary video-game programmer) surmised that some of the code might have been automatically translated from FORTRAN some years ago.

For example, on line 510 of SetupModel.cpp there is a loop over all the "places" the simulation knows about. This code appears to be trying to calculate R0 for "places". Hotels are excluded during this pass, without explanation.

This bit of code highlights an issue Caswell Bligh has discussed in your site's comments: R0 isn't a real characteristic of the virus. R0 is both an input to *and an output of* these models, and is routinely adjusted for different environments and situations. Models that consume their own outputs as inputs is problem well known to the private sector – it can lead to rapid divergence and incorrect prediction. There's a discussion of this problem in section 2.2 of the Google paper, "Machine learning: the high interest credit card of technical debt".

**Continuing development.** Despite being aware of the severe problems in their code that they "haven't had time" to fix, the Imperial team continue to add new features; for instance, the model attempts to simulate the impact of digital contact tracing apps.

Adding new features to a codebase with this many quality problems will just compound them and make them worse. If I saw this in a company I was consulting for I'd immediately advise them to halt new feature development until thorough regression testing was in place and code quality had been improved.

**Conclusions.** All papers based on this code should be retracted immediately. Imperial's modelling efforts should be reset with a new team that isn't under Professor Ferguson, and which has a commitment to replicable results with published code from day one.

On a personal level, I'd go further and suggest that all academic epidemiology be defunded. This sort of work is best done by the insurance sector. Insurers employ modellers and data scientists, but also employ managers whose job is to decide whether a model is accurate enough for real world usage and professional software engineers to ensure model software is properly tested, understandable and so on. Academic efforts don't have these people, and the results speak for themselves.

**My identity.** Sue Denim isn't a real person (read it out). I've chosen to remain anonymous partly because of the intense fighting that surrounds lockdown, but there's also a deeper reason. This situation has come about due to rampant credentialism and I'm tired of it. As the widespread dismay by programmers demonstrates, if anyone in SAGE or the Government had shown the code to a working software engineer they happened to know, alarm bells would have been rung immediately. Instead, the Government is dominated by academics who apparently felt unable to question anything done by a fellow professor. Meanwhile, average citizens like myself are told we should never question "expertise". Although I've proven my Google employment to Toby, this mentality is damaging and needs to end: please, evaluate the claims I've made for yourself, or ask a programmer you know and trust to evaluate them for you.

---

✉ **Subscribe** ▼                                                                                    Login

Please login to comment

**641 COMMENTS**                                                                  ⚡  🔥    Oldest ▼

**Will Jones**  🕐 9 months ago

Devastating. Heads must roll for this, and fundamental changes be made to the way government relates to academics and the standards expected of researchers. Imperial College should be ashamed of themselves.

**641**

👍 **409** 👎

**Lms2**  🕐 9 months ago

| 💬 *Reply to* *Will Jones*