

Bayesian model for COVID IFR

Witold Wiecek for 1 Day Sooner

Last updated 2020-09-03

Introduction

This is a short document describing a Bayesian model for synthesising information on many infection fatality rates (IFRs) into a single estimate that can be made specific to certain age groups or adjusted by co-morbidity status. The analysis presented here is a form of Bayesian meta-analysis, in that our primary objective is to weigh sources of evidence in a way that captures both variability (heterogeneity across different settings) and uncertainty.

Our ultimate objective is to characterise risks in a particular setting, population and time, in a way that is useful to understanding risks of human challenge trials (HCTs). Therefore, as a minimum, we want to incorporate variability into our prediction. Even better would be to understand how different factors can drive heterogeneity. Indeed, *a priori*, we can hypothesise that the three main drivers of differences in IFRs are time-specific, population-specific and otherwise country-specific.

The role of time may be due to new treatments, improvements over time in our ability to treat Covid-19 or selection pressures which may lead to more benign versions of the virus. Country-specific or location-specific factors in IFR data may be driven by under-reporting, health care factors (including access to health care services) or underlying distributions of known risk factors. Additionally, some unknown risk factors (e.g. genetic) may also be operating, in which case controlling for age and co-morbidities will be not sufficient to account for cross-location differences.

To address these drivers of differences in observed IFRs we develop a Bayesian model and apply it to publically available summary data on IFRs from multiple countries and contexts, with particular focus on the impact of age.

Methods

Bayesian model for evidence synthesis

What follows is an adaptation of typical methods of Bayesian evidence synthesis to analysis of IFRs. IFR is a proportion statistic, calculated as the ratio of deaths to infections in some population. Early estimates, e.g. by Verity et al. (2020), place it at over 0.6% globally. However, the risk of death is orders of magnitude higher in particular high risk groups, especially in the elderly, than in the general population.

We can use Bayesian models for repeated binary trials, accounting for the fact that different populations studies at different times have different average probability of events. We use hierarchical modelling framework to assume that the context-specific estimates of IFR_i (measured in different settings, with some uncertainty) are all linked using some common parameters.

The most straight-forward and “canonical” ways to implement such a Bayesian model is by modelling log odds of the event.¹ Deeks (2002) present a general treatment. Note, that for very rare events the odds of mortality are very similar to probability of mortality, but we model events on odds scale as a good “generic” approach to modelling binary data (in this case death following infections). Another advantage of such a model is that it can use either individual-level or summary data and work with covariates (such as gender, age, time of the study, co-morbidities), captured as odds ratios or risk ratios².

Basic models for analysis of binary data can be implemented using existing statistical analysis packages (see, for example, *baggr* by Wiecek and Meager (2020)), by treating IFR as a logit-normal parameter to meta-analyse. However, note that when no deaths are observed, analysis of the ratio statistic that is IFR (ratio of observed deaths to modelled infections) is problematic. Therefore we propose a “custom” model that built in Stan which treats deaths and *prevalences* as data (rather than the IFRs).

Let d_k denote observed deaths for data point k and assume that logit of prevalence p_k in the population of n_k subjects is obtained from some model. We can then write:

$$d_k \sim \text{Binomial}(n_k, p_k \text{IFR}_k) \quad (1)$$

$$\text{logit}(p_k) \sim \mathcal{N}(\mu_k^{(p)}, \sigma_k^{(p)}) \quad (2)$$

where $\sigma_k^{(p)}$ and $\mu_k^{(p)}$ are parameters derived from the existing models of prevalence.

The k data points collected can span many locations (studies); we denote them by loc_k and the total number of locations by N_{loc} . We can also collect other covariates impacting the IFRs, such as age groups (which we identify with median age of the population being studied, MedianAge_k). We denote all of the covariates using a design matrix X and denote by N_p the number of columns in X . We assume the impact on IFR is on logit scale, same as in the “canonical” logistic models of binary data that we mentioned above:

$$\text{logit}(\text{IFR}_k) = \theta_{\text{loc}_k} + X\beta$$

where θ is an N_{loc} -dimensional vector of location-specific (random) effects on IFR and β is N_p dimensional vector of (fixed) covariate effects.

We implement our model in Stan and assume very weakly informative priors on all parameters:

```
model {
  //Likelihood:
  // logit_ifr = theta_k[loc] + to_vector(X*beta);
  logit_prevalence ~ normal(mean_prevalence, sd_prevalence);
  obs_deaths ~ binomial(population, prevalence .* ifr);
  theta_k ~ normal(tau, sigma);

  //Priors:
  tau ~ normal(0, 10);
  sigma ~ normal(0, 10);
  beta ~ normal(0, 10);
}
```

¹It is also possible to work with IFR_i parameters and treat them as derived from Beta distribution with some “hyperparameters” α and β of Beta distribution, as done by e.g. Carpenter (2016). That approach, however, does not offer an easy way of modelling impact of covariates (e.g. age and co-morbidities) on the rates.

²If only summary data are available, covariates can be defined as study level distributions (e.g. % male)

Model data

We use estimates originally collected by Levin, Cochran, and Walsh (2020) to construct the first version of analysis dataset. The input data into our model consists of deaths (treated as known) and prevalences (treated as logit-distributed parameter with known mean and SD) in all reported age groups in all studies³.

```
## Inference for Stan model: ifr_with0.
## 4 chains, each with iter=2000; warmup=1000; thin=1;
## post-warmup draws per chain=1000, total post-warmup draws=4000.
##
##               mean se_mean   sd  2.5%  25%   50%   75% 97.5% n_eff Rhat
## tau         -8.66         0 0.20 -9.05 -8.79 -8.66 -8.53 -8.26 4392   1
## sigma        0.73         0 0.17  0.48  0.61  0.70  0.82  1.13 3767   1
## beta[1]      1.07         0 0.01  1.05  1.06  1.07  1.08  1.10 1939   1
##
## Samples were drawn using NUTS(diag_e) at Thu Sep 03 16:15:46 2020.
## For each parameter, n_eff is a crude measure of effective sample size,
## and Rhat is the potential scale reduction factor on split chains (at
## convergence, Rhat=1).
```

Predictive checks for the model

...

References

- Carpenter, Bob. 2016. “Hierarchical Partial Pooling for Repeated Binary Trials.” <https://mc-stan.org/users/documentation/case-studies/pool-binary-trials.html>.
- Deeks, Jonathan J. 2002. “Issues in the Selection of a Summary Statistic for Meta-Analysis of Clinical Trials with Binary Outcomes.” *Statistics in Medicine* 21 (11): 1575–1600. <https://doi.org/10.1002/sim.1188>.
- Levin, Andrew T., Kensington B. Cochran, and Seamus P. Walsh. 2020. “ASSESSING THE AGE SPECIFICITY OF INFECTION FATALITY RATES FOR COVID-19: META-ANALYSIS & PUBLIC POLICY IMPLICATIONS.” *medRxiv*, July, 2020.07.23.20160895. <https://doi.org/10.1101/2020.07.23.20160895>.
- Verity, Robert, Lucy C. Okell, Ilaria Dorigatti, Peter Winskill, Charles Whittaker, Natsuko Imai, Gina Cuomo-Dannenburg, et al. 2020. “Estimates of the Severity of Coronavirus Disease 2019: A Model-Based Analysis.” *The Lancet Infectious Diseases* 0 (0). [https://doi.org/10.1016/S1473-3099\(20\)30243-7](https://doi.org/10.1016/S1473-3099(20)30243-7).
- Wiecek, Witold, and Rachael Meager. 2020. “Baggr: Bayesian Aggregate Treatment Effects Package.” Zenodo. <https://doi.org/10.5281/zenodo.3813443>.

³This basic approach exaggerates uncertainty, as we treat different 95% intervals reported in the study as uncorrelated.

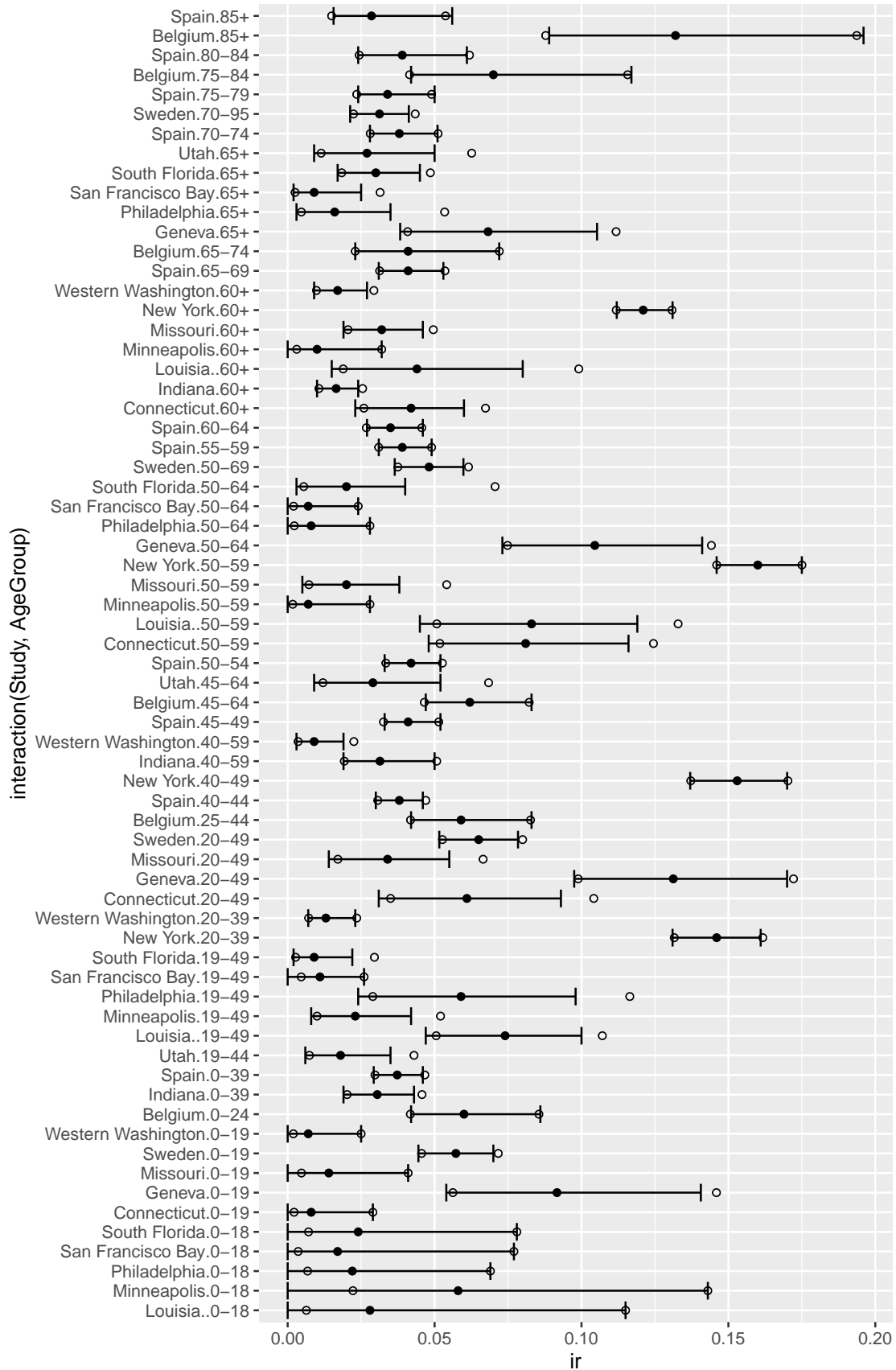


Figure 1: Distribution of model-estimated prevalences (95% CI's reported by modelling studies) collected by Levin, Cochran, and Walsh (2020). Additional points show 95% CIs recreated by assuming logit-normal distribution of prevalence.