Covid nowcasting

Jonas Wallin joint work with Adam Altmejd (SU), Joacim Rocklöv (UU)

> Home May 21, 2021



This talk is about nowcasting the number of deaths on a given with respect to what is reported so far.

- This talk is about nowcasting the number of deaths on a given with respect to what is reported so far.
- For a specific day it is only often after several days (weeks) that one has access to how many actually died that day.

- This talk is about nowcasting the number of deaths on a given with respect to what is reported so far.
- For a specific day it is only often after several days (weeks) that one has access to how many actually died that day.
- In Sweden and in the UK the data is presented before fully collected.

- This talk is about nowcasting the number of deaths on a given with respect to what is reported so far.
- For a specific day it is only often after several days (weeks) that one has access to how many actually died that day.
- In Sweden and in the UK the data is presented before fully collected.
- To have early access to the number of death is important for knowing where the epidemic is going and what efforts needs to be put in.

Antal avlidna per dag



Our goal is to model how many died on a given day. We denote:

 $dea_t-\ensuremath{\mathsf{actual}}$ number of indivuals that died of covid on day t

Our data is how many people are reported at a given day. We denote:

 $rep_{t,k}-$ How many indivuals reported dead on day t reported at day k

Thus let us examine how many died on 19th of May 2020 (t):

k:	Report date	$(\sum_{j=t}^{k} rep_{t,j})$
t	05 - 19	0
t+1	05 - 20	5
t+2	05 - 21	12
t+3	05 - 22	12
t+4	05 - 23	15
-	÷	
t + 349	2021 - 05 - 05	55



 Long time ago, I worked on estimating populations of natterjack toads, in Archipelago of Bohuslän.



- Long time ago, I worked on estimating populations of natterjack toads, in Archipelago of Bohuslän.
- Estimation of population is well studied problem in ecology. Typical methods are type of capture-recapture method.



- Long time ago, I worked on estimating populations of natterjack toads, in Archipelago of Bohuslän.
- Estimation of population is well studied problem in ecology. Typical methods are type of capture-recapture method.
- If p is probability of catching an animal and n is the number of catches and N the target population. Then approximately N ≈ n/p.



- Long time ago, I worked on estimating populations of natterjack toads, in Archipelago of Bohuslän.
- Estimation of population is well studied problem in ecology. Typical methods are type of capture-recapture method.
- If p is probability of catching an animal and n is the number of catches and N the target population. Then approximately $N \approx \frac{n}{p}$.
- On the islands there was small isolated populations. This meant that one could use rejection sampling or capture-retain method.

- One collects as many toads as one can (in a bucket). Then one waits a while and then one do it again. Repeat the procurer T times.
- This gives one a series of observation:

 $toad_1, toad_2, \ldots, toad_T$

- One collects as many toads as one can (in a bucket). Then one waits a while and then one do it again. Repeat the procurer T times.
- This gives one a series of observation:

 $toad_1, toad_2, \ldots, toad_T$

If one collects the toads from a small populations then one should get fewer toads each time and this actually gives some information of the actual size of the population.

- One collects as many toads as one can (in a bucket). Then one waits a while and then one do it again. Repeat the procurer T times.
- This gives one a series of observation:

 $toad_1, toad_2, \ldots, toad_T$

- If one collects the toads from a small populations then one should get fewer toads each time and this actually gives some information of the actual size of the population.
- Suppose the probability of catching a toad is p and the true population is N then:

 $toad_1 \sim Bin(N, p)$ $toad_2 \sim Bin(N - toad_1, p)$ $toad_3 \sim Bin(N - toad_1 - toad_2, p)$

Now we can estimate both p and N.

 \blacksquare In a Bayesian frame work we can estimate them using posterior distribution. Or treat \hat{N} and \hat{p} as parameters and estimate the ML

- In a Bayesian frame work we can estimate them using posterior distribution. Or treat \hat{N} and \hat{p} as parameters and estimate the ML
- I realised yesterday that the ML is clearly biased here but it seems easy to correct this, but that for an other day maybe (Does not matter for our data since number of observations will be large).
- Let for example be Bayesian and study an example.

$$toad = [22, 15]$$

posterior:



0.014

0:005 bit

100

Ν

150

200

$$toad = [22, 15, 14]$$





$$toad = [22, 15, 14, 9]$$



$$\pi(N, p|toad) \propto \left(\prod_{i}^{m} \binom{N - \sum_{1}^{i-1} toad_{i}}{toad_{i}}\right)_{\substack{a \in \mathcal{B} \\ p \sum_{i} toad_{i}}(1-p)^{mN - \sum_{i} toad_{i}}} \left(\prod_{i=1}^{\mathfrak{B}} \frac{1}{p^{i-1}}\right)_{\substack{a \in \mathcal{B} \\ 0,1 \dots 0,2 \dots 0,3 \dots 0,4 \dots 0,5 \dots 0,4 \dots 0,5 \dots 0,4 \dots 0,5 \dots 0,5$$

$$toad = [22, 15, 14, 9, 7]$$



$$\pi(N, p|toad) \propto \left(\prod_{i}^{m} \left(N - \sum_{1}^{i-1} toad_{i}\right)\right)_{\text{a}} \int_{g}^{g} \int_{0,1}^{g} \int$$

$$toad = [22, 15, 14, 9, 7, 6]$$



$$\pi(N, p|toad) \propto \left(\prod_{i}^{m} \binom{N - \sum_{1}^{i-1} toad_{i}}{toad_{i}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} \frac{1}{p^{\sum_{i} toad_{i}} (1-p)^{mN - \sum_{i} toad_{i}}}\right)_{\mathbb{R}} \left(\sum_{g \in \mathcal{F}_{i}} \frac{1}{p^{\sum_{i} toad_{i}} \frac{1}{p^$$

Back to the covid-data

 For each given day we can apply the same idea for the reporting

$$rep_{t,\cdot} = [rep_{t,t}, rep_{t,t+1}, \dots, rep_{t,T}].$$

and thus estimate dea_t using the model

÷

$$rep_{t,t} \sim Bin(dea_t, p_0)$$

 $rep_{t,t+1} \sim Bin(dea_t - rep_{t,t}, p_1)$

$$rep_{t,t+i} \sim Bin(dea_t - \sum_{j=0}^{i-1} rep_{t,t+j}, p_i)$$

Back to the covid-data

 For each given day we can apply the same idea for the reporting

$$rep_{t,\cdot} = [rep_{t,t}, rep_{t,t+1}, \dots, rep_{t,T}].$$

and thus estimate dea_t using the model

÷

$$rep_{t,t} \sim Bin(dea_t, p_0)$$

 $rep_{t,t+1} \sim Bin(dea_t - rep_{t,t}, p_1)$

$$rep_{t,t+i} \sim Bin(dea_t - \sum_{j=0}^{i-1} rep_{t,t+j}, p_i)$$

Back to the covid-data

For each given day we can apply the same idea for the reporting

$$rep_{t,\cdot} = [rep_{t,t}, rep_{t,t+1}, \dots, rep_{t,T}].$$

and thus estimate dea_t using the model

÷

$$rep_{t,t} \sim Bin(dea_t, p_0)$$
$$rep_{t,t+1} \sim Bin(dea_t - rep_{t,t}, p_1)$$

$$rep_{t,t+i} \sim Bin(dea_t - \sum_{j=0}^{i-1} rep_{t,t+j}, p_i)$$

I had missed this for medical data this had been done in previously in [Lawless, J. F. 1994]. (Of course they missed the biological connection dating back to [Moran, P. A. P. 1951])

 Often the data when working with Binomial distribution is overdispersed (higher variability compared to the distribution), here there is no exception.

- Often the data when working with Binomial distribution is overdispersed (higher variability compared to the distribution), here there is no exception.
- Once enough data is collected one can assume that one have the truth. In uk we are very stable after 30 days (much sooner). In this case things reduces to that both rep_t, and dea_t is known and we only need to estimate p.

- Often the data when working with Binomial distribution is overdispersed (higher variability compared to the distribution), here there is no exception.
- Once enough data is collected one can assume that one have the truth. In uk we are very stable after 30 days (much sooner). In this case things reduces to that both rep_t, and dea_t is known and we only need to estimate p.
- Let us examine how many are reported at lag 2 (number of cases reported two days after).

- Often the data when working with Binomial distribution is overdispersed (higher variability compared to the distribution), here there is no exception.
- Once enough data is collected one can assume that one have the truth. In uk we are very stable after 30 days (much sooner). In this case things reduces to that both rep_t, and dea_t is known and we only need to estimate p.
- Let us examine how many are reported at lag 2 (number of cases reported two days after).
- With a bit of simplification we approximately have:

$$\hat{p}_2 = \frac{\sum_t rep_{t,t+2}}{\sum_t dea_t} \approx \frac{\sum_t rep_{t,t+2}}{\sum_t \sum_{i=0}^{30} dea_{t,t+i}}$$

 \blacksquare Given \hat{p} we have approximately

$$rep_{t,t+2} \sim Bin(rep_{t,t+30} \approx dea_t, \hat{p}_2)$$

Thus we generate CI for $\frac{rep_{t,t+2}}{dea_t}$ to check the Binomial distribution assumption.

 \blacksquare Given \hat{p} we have approximately

$$rep_{t,t+2} \sim Bin(rep_{t,t+30} \approx dea_t, \hat{p}_2)$$

Thus we generate CI for $\frac{rep_{t,t+2}}{dea_t}$ to check the Binomial distribution assumption.



dates

To address this issue we instead of assuming a Binomial distribution use a Beta-Binomial distribution.

- To address this issue we instead of assuming a Binomial distribution use a Beta-Binomial distribution.
- This corresponds to:

$$p \sim Beta(\mu, M),$$
$$rep_{t,t+2} \sim Bin(\sum_{i=0}^{30} dea_{t,t+i}rep_{t,t+i} \approx dea_t, p).$$

One can integrate out \boldsymbol{p} and get

$$rep_{t,t+2} \sim BB(\sum_{i=0}^{30} dea_{t,t+i}rep_{t,t+i} \approx dea_t, \mu, M).$$

here μ is the expected value while M is overdispersion.

- To address this issue we instead of assuming a Binomial distribution use a Beta-Binomial distribution.
- This corresponds to:

$$p \sim Beta(\mu, M),$$

$$rep_{t,t+2} \sim Bin(\sum_{i=0}^{30} dea_{t,t+i} rep_{t,t+i} \approx dea_t, p).$$

One can integrate out \boldsymbol{p} and get

$$rep_{t,t+2} \sim BB(\sum_{i=0}^{30} dea_{t,t+i} rep_{t,t+i} \approx dea_t, \mu, M).$$

here μ is the expected value while M is overdispersion.





Sweden

For Sweden things are a bit more complicated due to clear time trend in the data (at the beginning of the pandemic)



dates

Sweden

- For Sweden things are a bit more complicated due to clear time trend in the data (at the beginning of the pandemic)
- We will ignore this for now and focus on the second wave.



dates

• So our model at day t is now

$$rep_{t} \sim BB(dea_t, \theta_{BB})$$

where θ_{BB} is set of parameters controlling for different lags, report days etc.

• Let us examine the result for Sweden at 17th of November 2020. Where we fitted θ_{BB} with ML removing the thirty latest dates and then sample from the posterior distribution i.e.

$$dea_t | rep_{t,}, \theta_{BB} \sim BB(rep_{t,}; dea_t, \theta_{BB})$$

result



date

It is easy to see that the deaths are correlated which the model does not take into account.

result

- It is easy to see that the deaths are correlated which the model does not take into account.
- The underlying factor is the number of infected (also by age group matters).

result

- It is easy to see that the deaths are correlated which the model does not take into account.
- The underlying factor is the number of infected (also by age group matters).
- We can separately (from the reporting) try to model the epidemic.

$$\frac{dea_t|dea_{t-1:1} \sim ?}{rep_{t,} \sim BB(dea_t, \theta_{BB})}$$

Adding other information

- We also have other data that correlates with the deaths: ICU, hospitalized with Covid, and infected in eldercare.
- Let M_t denote the one of these time series (or a smoothed version) we can use

$$dea_t \sim Po(exp\left(\beta_0 + \beta_1 M_{t-c}\right))$$

here we need to estimate β_0, β_1, c .

- We in-fact included $M_t = ICU_t$ in our model.... at the completely wrong time.
- Issue they only correlate with the dea_t and this correlation can change over time

Adding other information



We instead would like to incorporate an underlying pandemic trend.

- We instead would like to incorporate an underlying pandemic trend.
- Ideally we would like to model as follows

$$\begin{split} X(t) &\sim GP(t; \theta_{GP}), \\ \lambda_t &= \exp(X(t)), \\ dea_t &\sim Poisson(\lambda_t), \\ rep_{t,} &\sim BB(dea_t, \theta_{BB}). \end{split}$$

- We instead would like to incorporate an underlying pandemic trend.
- Ideally we would like to model as follows

$$\begin{split} X(t) &\sim GP(t; \theta_{GP}), \\ \lambda_t &= \exp(X(t)), \\ dea_t &\sim Poisson(\lambda_t), \\ rep_{t,} &\sim BB(dea_t, \theta_{BB}). \end{split}$$

■ Then sampling from the posterior distribution $dea_t | rep_t$, one gets predictions.

- We instead would like to incorporate an underlying pandemic trend.
- Ideally we would like to model as follows

$$\begin{split} X(t) &\sim GP(t; \theta_{GP}), \\ \lambda_t &= \exp(X(t)), \\ dea_t &\sim Poisson(\lambda_t), \\ rep_{t,} &\sim BB(dea_t, \theta_{BB}). \end{split}$$

- Then sampling from the posterior distribution $dea_t | rep_t$, one gets predictions.
- However, this requires a bit of coding and also mentoring of the MCMC chains which we (I) am to lazy to do.

- We instead would like to incorporate an underlying pandemic trend.
- Ideally we would like to model as follows

$$\begin{split} X(t) &\sim GP(t; \theta_{GP}), \\ \lambda_t &= \exp(X(t)), \\ dea_t &\sim Poisson(\lambda_t), \\ rep_{t,} &\sim BB(dea_t, \theta_{BB}). \end{split}$$

- Then sampling from the posterior distribution $dea_t | rep_t$, one gets predictions.
- However, this requires a bit of coding and also mentoring of the MCMC chains which we (I) am to lazy to do.
- Instead we look to approximate steps in the chain to get a faster algorithm.

 Since the data seems to be Poisson (or negative binomial) distributed we use a square root transformation

$$dea_t^{sqrt} = \sqrt{dea_t}$$

this since it gives the right relation between variance and mean.

Since the data seems to be Poisson (or negative binomial) distributed we use a square root transformation

$$dea_t^{sqrt} = \sqrt{dea_t}$$

this since it gives the right relation between variance and mean.Then the latent part becomes

$$\begin{aligned} X(t) &\sim GP(t;\theta) \\ dea_t^{sqrt} &\sim \mathcal{N}\left(X(t),\sigma^2\right) \end{aligned}$$

Since the data seems to be Poisson (or negative binomial) distributed we use a square root transformation

$$dea_t^{sqrt} = \sqrt{dea_t}$$

this since it gives the right relation between variance and mean.Then the latent part becomes

$$X(t) \sim GP(t; \theta)$$
$$dea_t^{sqrt} \sim \mathcal{N}\left(X(t), \sigma^2\right)$$

Yet the likelihood component

$$rep_{t}$$
, ~ $BB(dea_t, \theta_{BB})$

is also still intractable and this we will approximate in the next step.

Here it is easy to sample from $dea_t | rep_{t_i}$ from the likelihood only model (denoted \mathcal{F}_{simple}):

$$rep_{t,} \sim BB(dea_t, \theta_{BB}).$$

From these samples we approximate the posterior distribution with the following

$$f(\sqrt{dea_t}|rep_{t,}) \approx \mathcal{N}\left(\sqrt{dea_t}; \mu_t, \sigma^2\right)$$

Where

$$\mu_t = \mathbb{E}^{MC} \left[\sqrt{dea}_t | rep_{t,}, \mathcal{F}_{simple} \right]$$

and

$$\sigma_t^2 = \mathbb{V}^{MC} \left[\sqrt{dea}_t | rep_{t,}, \mathcal{F}_{simple} \right]$$

























further issues

We are using a metric for prediction fit known as continuous probability rank score (CRPS).



date

Sweden model final



Source: Folkhälsomyndigheten and ECDC, Updated: 2021-05-17, Latest version available at https://adamaltmeid.se/covid.

Uk data final



Source: GOV.UK. Updated: 2021-05-17.