

# State estimation of Covid-19 disease in Denmark

Kasper Schou Telkamp

2022-12-23

## Contents

<b>Introduction</b>	<b>2</b>
<b>Data</b>	<b>2</b>
<b>Methodology</b>	<b>4</b>
Modelling . . . . .	4
Parameter estimation . . . . .	6
<b>Results</b>	<b>6</b>
Estimating the total number of new hospital admissions in Denmark . . . . .	6
Estimating the number of new hospital admissions in Denmark grouped by region . . . . .	7
Residual analysis . . . . .	9
<b>Discussion</b>	<b>11</b>
<b>Conclusion</b>	<b>11</b>
<b>References</b>	<b>12</b>
<b>Appendix A</b>	<b>12</b>
<b>Appendix B</b>	<b>13</b>
<b>Appendix C</b>	<b>14</b>

# Introduction

Early detection of outbreaks with communicable diseases are of great importance in order to initiate timely interventions and help prevent disease spread. When dealing with vast amount of data, automated procedures can supplement traditional surveillance methods. In this report non-normal mixed effects models is evaluated on their ability to identify outbreaks of Covid-19 disease using data over new hospital admissions with Covid-19 in Denmark. Different implementations of generalized linear mixed models (GLMMs) in R packages is compared. Namely, the `glmmTMB` (Brooks et al., 2017) and `KFAS` (Helske, 2017) R packages available at Comprehensive R Archive Network (CRAN).

## Data

In this project, the daily record of new hospital admissions with Covid-19 in Denmark grouped by region of residence and totals are used. This report is based on data from 2020-03-01 to 2022-11-06 is used. The head and tail of the processed data are listed in Table 1.

Table 1: Processed data containing the daily record of new hospital admissions with Covid-19 in Denmark grouped by region of residence and with totals.

Dato	Hovedstaden	Sjælland	Syddanmark	Midtjylland	Nordjylland	Ukendt.Region	Total
2020-03-01	1	0	0	0	0	0	1
2020-03-02	0	0	0	0	0	0	0
2020-03-03	1	0	0	0	0	0	1
2020-03-04	0	0	0	0	0	0	0
...	...	...	...	...	...	...	...
2022-11-03	20	12	9	10	2	1	54
2022-11-04	13	15	6	3	3	0	40
2022-11-05	7	6	6	3	3	0	25
2022-11-06	8	10	7	4	2	0	31

The data is publicly available and were obtained from Statens Serum Institut (SSI) website<sup>1</sup>. SSI collects the data from the National Patient Registry (NPR), which contains information about outpatient contacts from Danish public as well as private hospitals. The data from NPR has some delay. Therefore, the inventory is updated daily with real-time data from the regions. The regions provide snapshot-data twice to SSI daily at 7am and 3pm. A hospital admission related to Covid-19 is defined as an admission, where a patient is admitted within 14 days after a positive SARS-CoV-2 test. Patients that are tested positive for SARS-CoV-2 during an admission is also registered as a Covid-19 related admission. Furthermore, admissions with Covid-19 are only registered for patients that are present in at least one snapshot, or if the patient have been admitted for more than 12 hours according to NPR. The total number of new admissions to the hospital in Denmark are visualized in Figure 1.

<sup>1</sup><https://covid19.ssi.dk/overvagningsdata/download-fil-med-overvaegningdata>

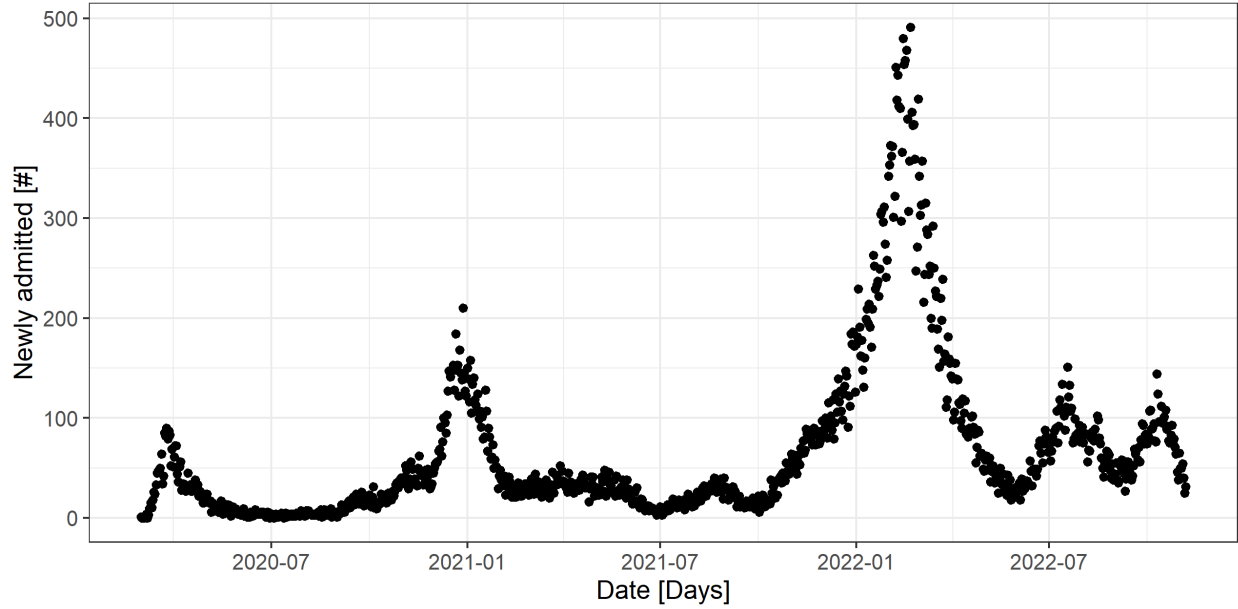


Figure 1: Total number of new hospital admissions in Denmark.

Clearly, the observed number of new hospital admissions are correlated in time. In the past consecutive epidemiological waves of Covid-19 disease have been observed, but with varying peak number of new hospital admissions. The peak number of new hospital admissions depends upon a number of influencing factors such as:

- Severity of illness
- Population in which the virus spreads
- How many are infected

Whereas the first two waves of new admissions to hospitals with Covid-19 disease can be explained by the severity of the illness and the population in which the virus spread, the large wave observed in early January, 2022, is more like likely caused by the vast amount of disease spread and how many that were infected at the same time.

Denmark is grouped into five regions: *Hovestaden*, *Midtjylland*, *Nordjylland*, *Sjælland*, and *Syddanmark*. If a patient is admitted to the hospital with Covid-19, but does not reside in any of the regions, they are marked as *Ukendt Region*. In Figure 2 the daily number of new hospital admissions grouped by region of residence is visualized.

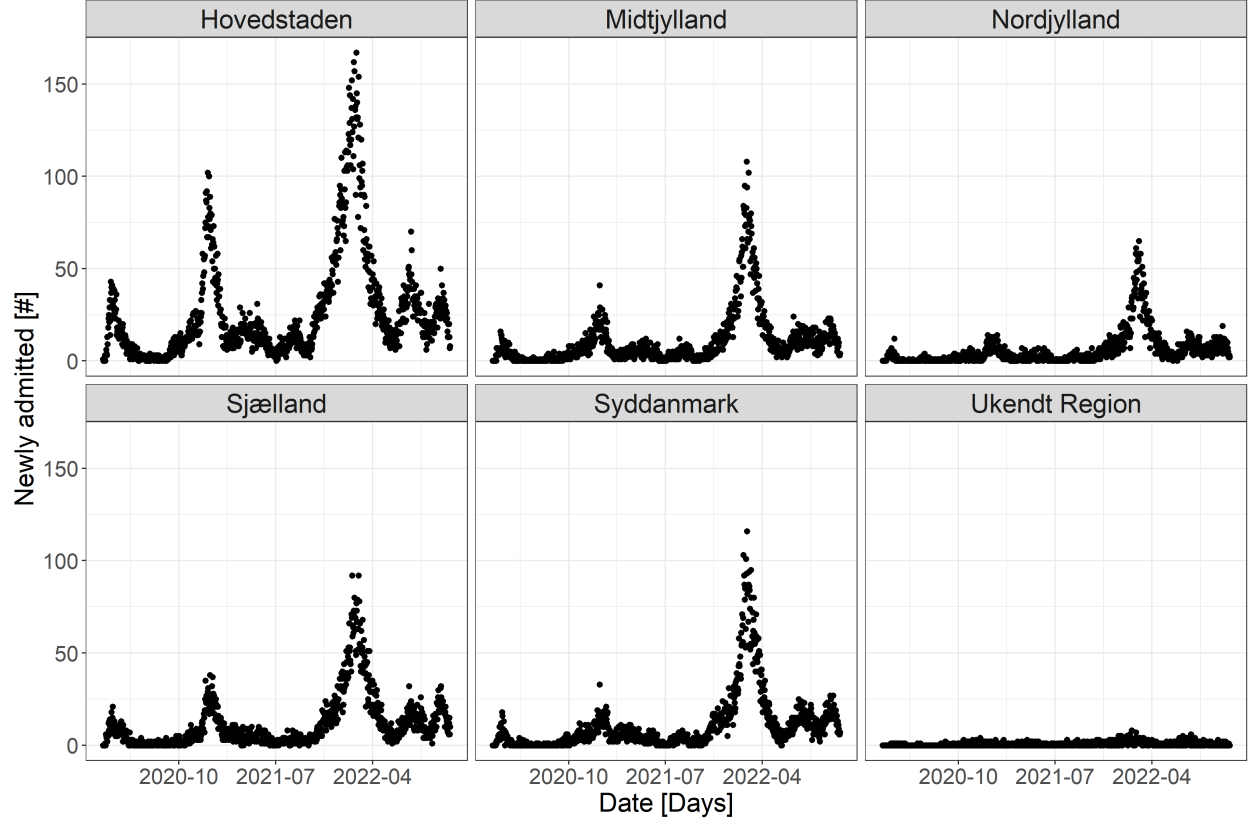


Figure 2: Total number of new hospital admissions in Denmark grouped by region of residence.

It can be seen that most of the new hospital admissions can be linked to *Hovedstaden*, while only a negligible amount are linked to the *Ukendt Region*. Intermediate amount of new hospital admissions are linked to *Midtjylland*, *Nordjylland*, *Sjælland*, and *Syddanmark*. This is largely due to the fact, that *Hovedstaden* considers more individuals and that the exposure therefore is higher in this region. Additionally, a correlation between the regions are observed.

## Methodology

In this section the GLMMs for modelling the number of new hospital admissions are formulated. Moreover, methods for approximating the likelihood functions and implementations in R packages for parameter estimation are presented.

### Modelling

In order to analyze the data a simple state space model is proposed. The count observations  $h_t$ ,  $t = 1, \dots, n$  in a period of  $n = 981$  days starting on the 1st of March, 2020, of new hospital admissions is assumed to follow a Poisson distribution  $h_t \sim P(\lambda_t)$  with intensities given by

$$\log(\lambda_t) = \beta + u_t \quad (1)$$

Here  $\beta$  is a fixed parameter, that represents the average intensity and  $u_t$  is a random effect that is assumed to follow a first order auto-regressive process

$$u_t = au_{t-1} + \epsilon_t \quad (2)$$

where  $\epsilon \sim N(0, \sigma^2)$ ,  $t > 1$  is a white noise process, and  $a$  and  $\sigma$  are model parameters. Using the results from Madsen (2007), it is assumed that  $u_1$  follow the stationary distribution of the first order auto-regressive process  $u_1 \sim N(0, \sigma^2/(1 - a^2))$ . Hence the joint likelihood becomes

$$L(a, \beta, \sigma; \mathbf{u}, \mathbf{y}) = \phi_{0, \frac{\sigma^2}{1-a^2}}(u_1) \prod_{t=2}^n (\phi_{0, \sigma^2}(u_t - au_{t-1})) \prod_{t=1}^n (p_{\lambda_t}(h_t)) \quad (3)$$

where  $\phi_{\mu, \sigma^2}$  is the probability density function (pdf) of the normal distribution with mean  $\mu$  and variance  $\sigma^2$ , and  $p_{\lambda}$  is the pdf of the Poisson distribution with mean  $\lambda$ . Intuitively, the model can be extended by modelling the individual regions. Hence,  $(a, \beta, \sigma)$  are 6-dimensional vectors.

In order to obtain the likelihood for the model parameters  $(a, \beta, \sigma)$  the observed random effects are integrated out. Hence, the marginal likelihood is obtained

$$L_M = (a, \beta, \sigma; \mathbf{y}) = \int_{\mathbb{R}^q} L(a, \beta, \sigma; \mathbf{u}, \mathbf{y}) d\mathbf{u} \quad (4)$$

where  $q$  is the number of random effects and  $a$ ,  $\beta$ , and  $\sigma$  are the parameters to be estimated.

In order to make computation of the joint likelihood function in (4) feasible, the estimation is carried out using the multivariate Laplace approximation.

### Laplace approximation

The marginal log-likelihood  $l_M(a, \beta, \sigma; \mathbf{y}) = \log(L_M(a, \beta, \sigma; \mathbf{y}))$  is approximated by a second order Taylor approximation around the optimum  $\tilde{\mathbf{u}} = \hat{\mathbf{u}}_{\theta}$  of the log-likelihood function w.r.t. the unobserved random variables  $\mathbf{u}$ , i.e.,

$$l(a, \beta, \sigma; \mathbf{u}, \mathbf{y}) \approx l(a, \beta, \sigma; \tilde{\mathbf{u}}, \mathbf{y}) - \frac{1}{2}(\mathbf{u} - \tilde{\mathbf{u}})^T H(\tilde{\mathbf{u}})(\mathbf{u} - \tilde{\mathbf{u}}) \quad (5)$$

where the first-order term of the Taylor expansion disappears since the expansion is done around the optimum  $\tilde{\mathbf{u}}$  and  $H(\tilde{\mathbf{u}}) = -l''_{uu}(a, \beta, \sigma, \mathbf{u}, \mathbf{y})|_{\mathbf{u}=\tilde{\mathbf{u}}}$  is the negative Hessian of the joint log-likelihood evaluated at  $\tilde{\mathbf{u}}$ .

Using the approximation in (5) on (4) the Laplace approximation of the marginal log-likelihood becomes (See Madsen & Thyregod (2011))

$$l_{M,LA}(a, \beta, \sigma; \mathbf{y}) = \log \int_{\mathbb{R}^q} \exp \left( l(a, \beta, \sigma; \tilde{\mathbf{u}}, \mathbf{y}) - \frac{1}{2}(\mathbf{u} - \tilde{\mathbf{u}})^T H(\tilde{\mathbf{u}})(\mathbf{u} - \tilde{\mathbf{u}}) \right) d\mathbf{u} \quad (6)$$

$$= l(a, \beta, \sigma; \tilde{\mathbf{u}}, \mathbf{y}) - \frac{1}{2} \log \left| \frac{H(\tilde{\mathbf{u}})}{2\pi} \right| \quad (7)$$

### Importance sampling

Importance sampling is a re-weighting technique for approximating integrals w.r.t. a density  $f$  by simulation in cases where it is not feasible to simulate from the distribution with density  $f$ . Instead it uses samples from a different distribution with density  $g$ , where the support of  $g$  includes the support of  $f$ .

## Parameter estimation

In this section two R packages for used to estimate the parameters are presented. Namely, **glmmTMB** and **KFAS** which are available at CRAN.

### glmmTMB

This section describes the R package **glmmTMB** by Brooks et al. (2017) for linear and GLMMs using Template Model Builder (TMB). The models are estimated using maximum likelihood estimation via TMB. Random effects are assumed to be Gaussian on the scale of the linear predictor and are integrated out using Laplace approximation. Additionally, gradients are calculated using automatic differentiation.

### KFAS

This section goes into detail with the R package **KFAS** by Helske (2017) for state space modelling with observations from the exponential family. The **KFAS** package can perform Kalman filtering and smoothing with exact diffuse initialization using an univariate approach.

In **KFAS** the Poisson distribution with intensity  $\lambda_t$  and exposure term  $e_t$  together with the log-link is supported. Thus we have  $E(h_t | \log(\lambda_t)) = \text{Var}(h_t | \log(\lambda_t)) = e_t \lambda_t$ . In this report the exposure term is assumed to be constant, i.e.  $e_t = 1$ . Hence, the differences are represented directly in the estimated parameters and latent state.

In order to make inferences of the Poisson model, **KFAS** finds a Gaussian model with the same conditional posterior mode as  $P(\lambda | \mathbf{h})$ . This is done through an iterative process with Laplace approximation of  $P(\lambda | \mathbf{h})$ , where the updated estimates of  $\log(\lambda_t)$  are computed via the Kalman filtering and smoothing from the approximating Gaussian model. The final estimates of  $\log(\hat{\lambda}_t)$  correspond to the mode of  $P(\lambda | \mathbf{h})$ . Generally, the difference between the mode and the mean is negligible. Nevertheless, our interest is focused on the intensity  $\lambda_t$  rather than the linear predictor  $\log(\lambda_t)$ .

Direct transformation from the linear predictor to the intensity introduces some bias. To solve this problem **KFAS** also contains methods based on importance sampling.

## Results

In this section the obtained parameter estimates are presented.

### Estimating the total number of new hospital admissions in Denmark

Table 2: That and this

Parameter	glmmTMB	KFAS
$\beta$	3.175 (0.726)	3.605 (0.02)

The smoothed estimate of the latent state estimated from **glmmTMB** and **KFAS** is visualized in Figure 3.

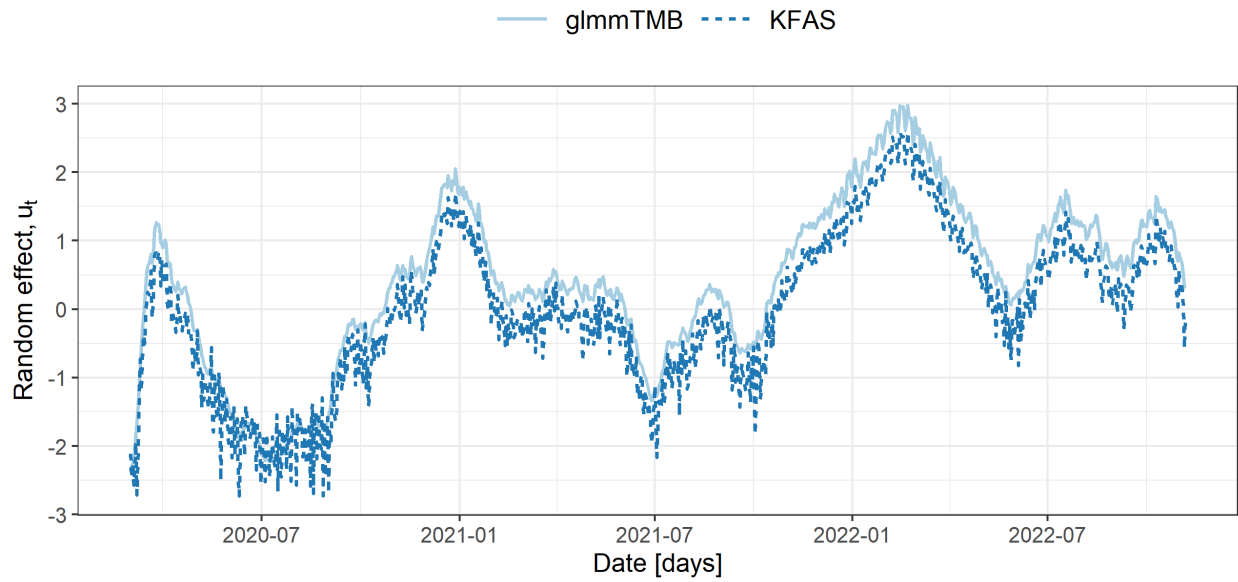


Figure 3: Smoothed estimates of the latent states from `glmmTMB` and `KFAS`.

In Figure 4 the total number of new hospitals admissions in Denmark is visualized together with the smoothed estimates from `glmmTMB` and `KFAS`.

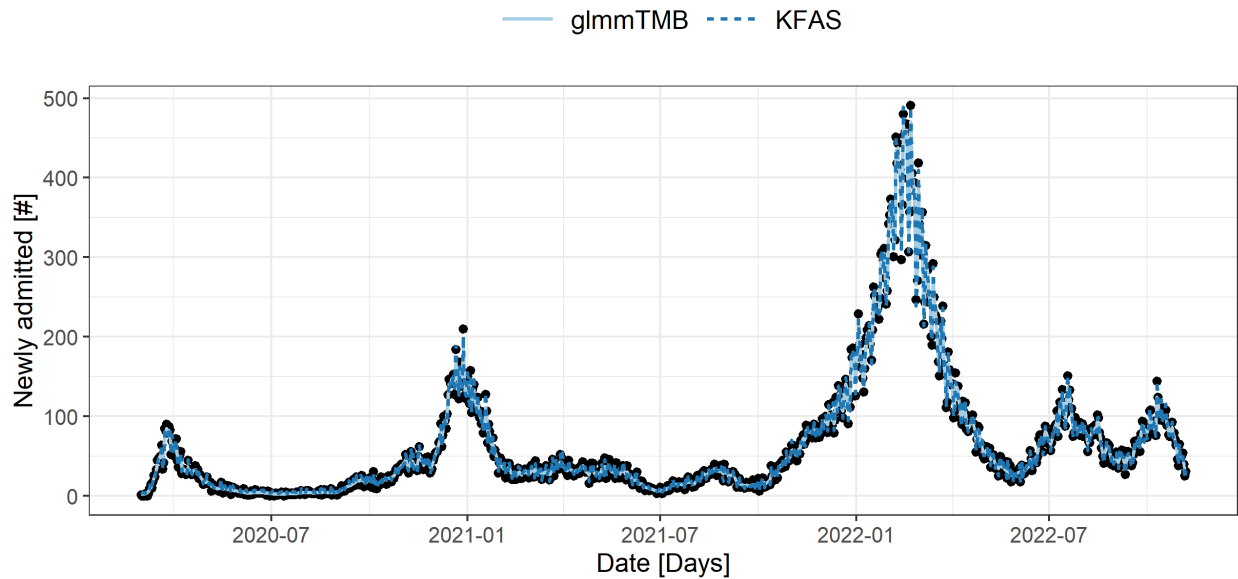


Figure 4: Total number of new hospital admissions in Denmark with smoothed estimates from `glmmTMB` and `KFAS`.

**Estimating the number of new hospital admissions in Denmark grouped by region**

Table 3: This and that

Parameter	glmmTMB	KFAS
$\beta_1$	2.438 (0.612)	2.625 (0.025)
$\beta_2$	1.428 (0.611)	1.583 (0.026)
$\beta_3$	0.818 (0.612)	0.688 (0.089)
$\beta_4$	1.65 (0.61)	1.805 (0.02)
$\beta_5$	1.354 (0.611)	1.503 (0.034)
$\beta_6$	-0.971 (0.619)	-0.848 (0.055)

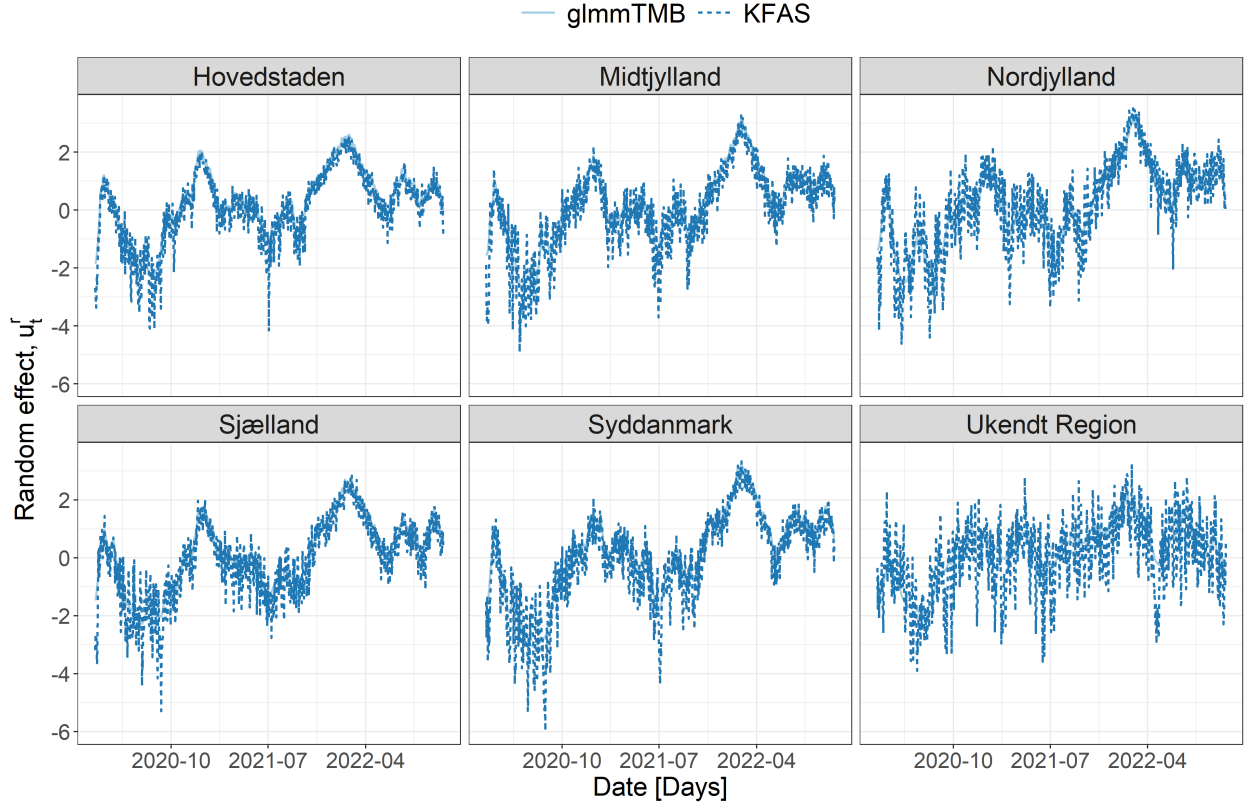


Figure 5: Smoothed estimates of the latent states from `glmmTMB` and `KFAS`.

In Figure 6 the total number of new hospitals admissions in Denmark grouped by region is visualized together with the smoothed estimates from `glmmTMB` and `KFAS`.



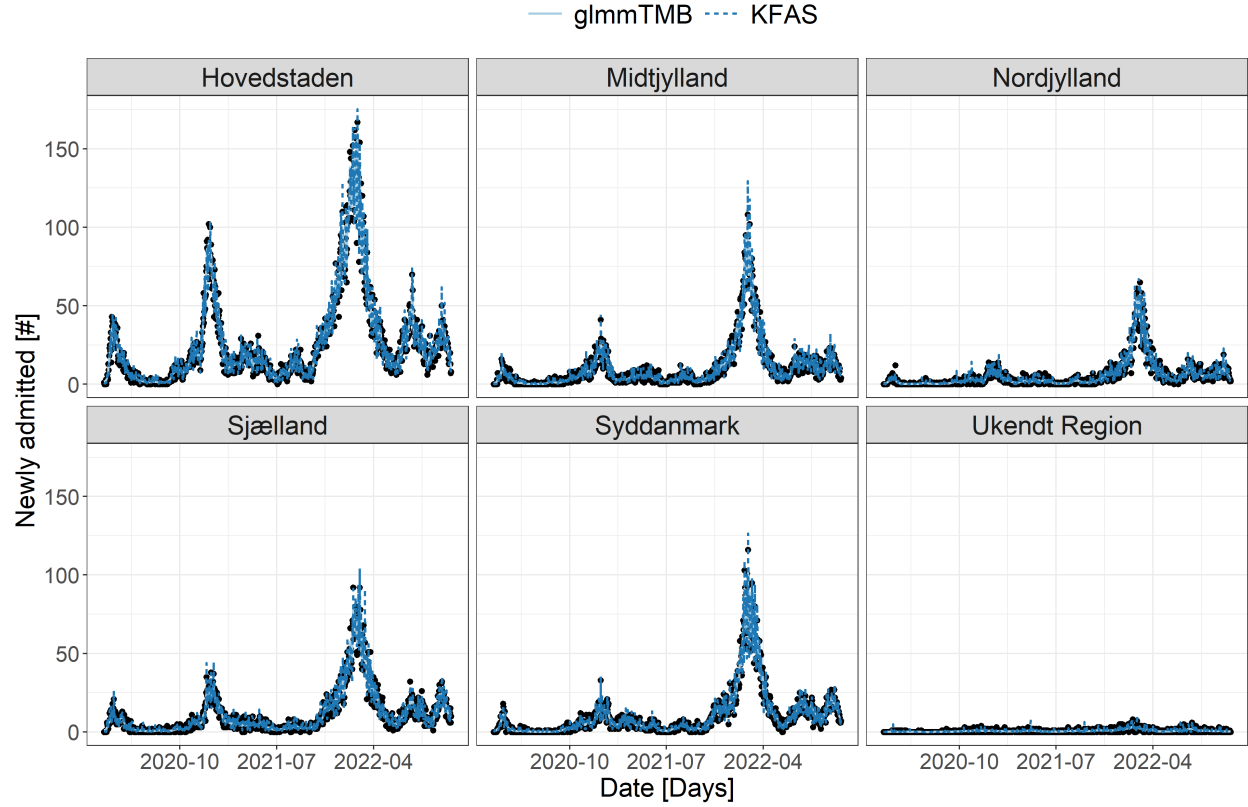


Figure 6: Total number of new hospital admissions in Denmark grouped by region with smoothed estimates from glmmTMB and KFAS.

## Residual analysis

Diagnostic plots of the one-step prediction residuals for the total number of new hospital admissions in Denmark are visualized in Figure 7. Overall the plots hint that the model can be improved. In A) and B) a significant auto-correlation in the residuals are detected. This is likely caused by the increased variance in the number of new admissions during the Covid-19 waves. The two lower plots, C) and D), indicate heavy tailed residuals.

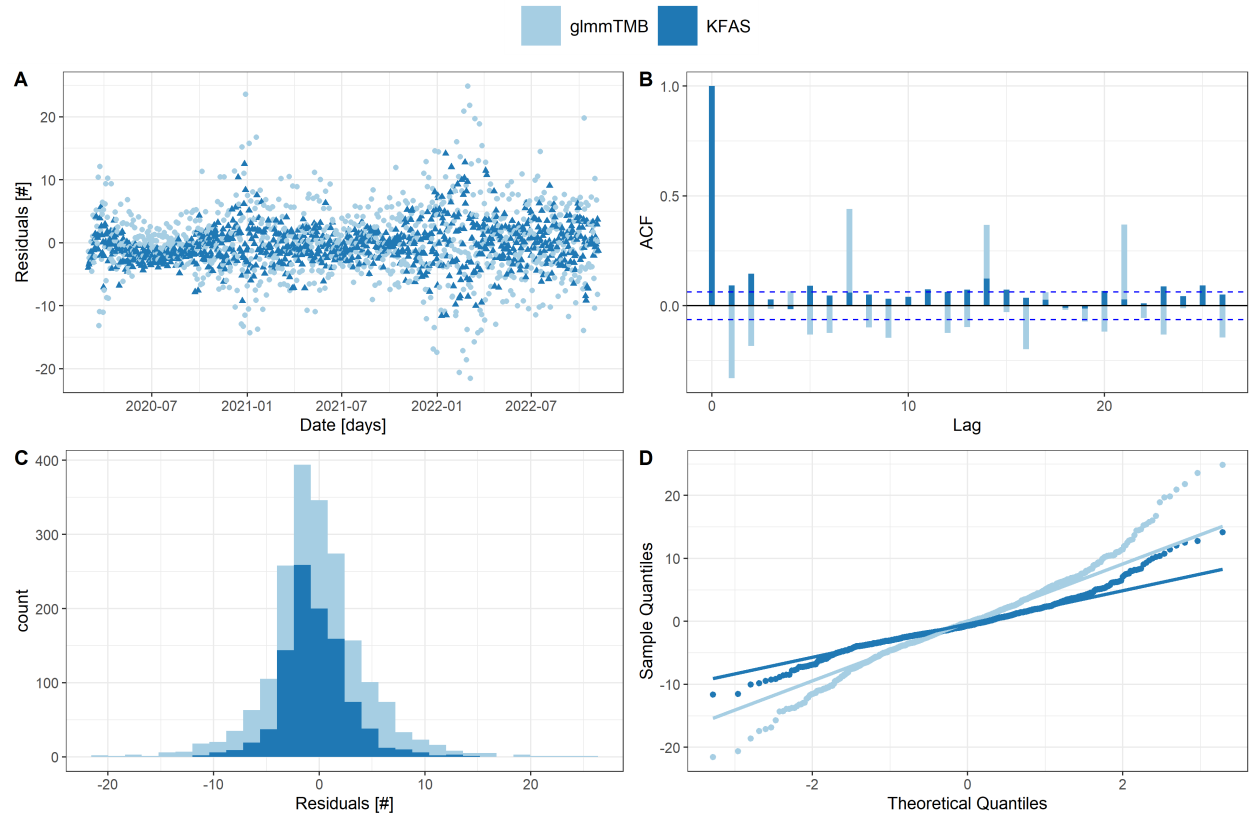


Figure 7: Diagnostic plots of the in-sample one-step prediction residuals. A) Time series of the residuals, B) auto-correlation function of the residuals, C) Histogram of the residuals, and D) Qunatile-Quantile plot of the residuals.

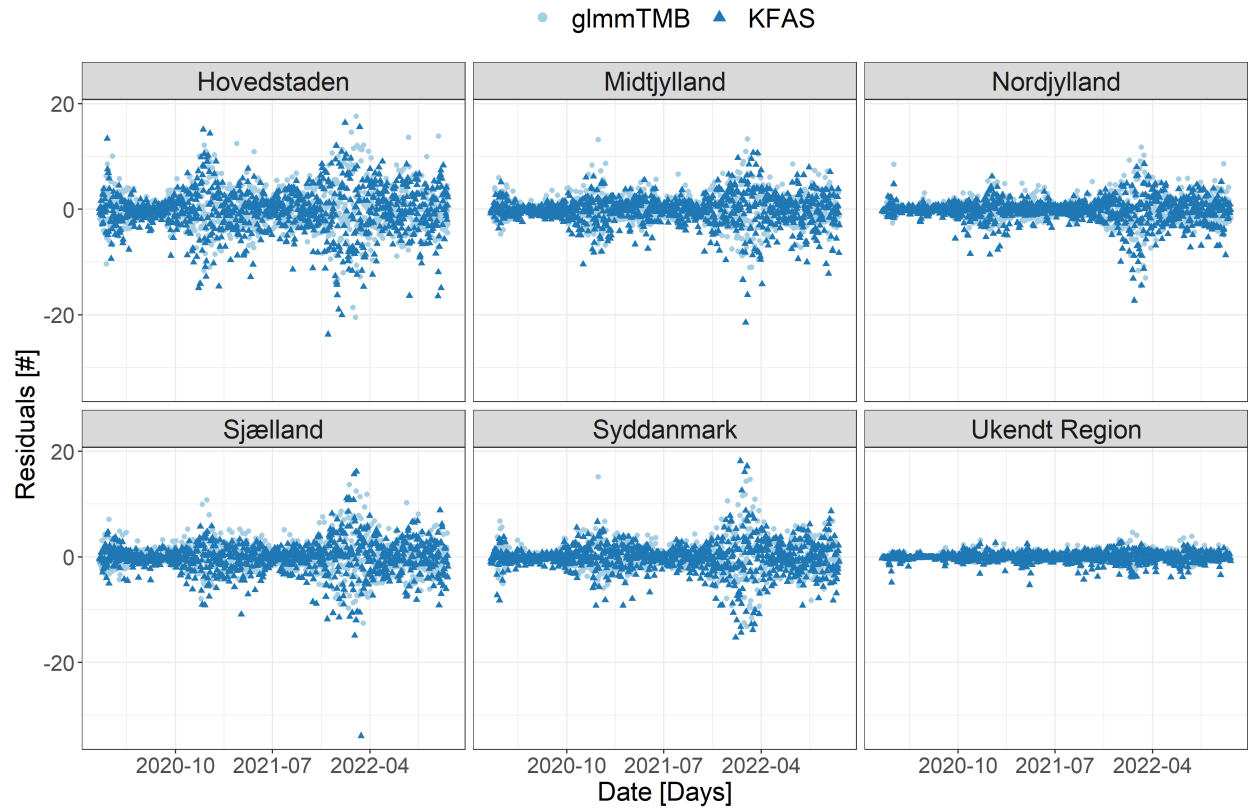


Figure 8: Model residuals from `glmmTMB` and `KFAS`.

## Discussion

Write it later

## Conclusion

## References

- Brooks, M. E., Kristensen, K., van Benthem, K. J., Magnusson, A., Berg, C. W., Nielsen, A., Skaug, H. J., Maechler, M., & Bolker, B. M. (2017). glmmTMB balances speed and flexibility among packages for zero-inflated generalized linear mixed modeling. *The R Journal*, 9(2), 378–400. <https://doi.org/10.32614/RJ-2017-066>
- Helske, J. (2017). KFAS: Exponential family state space models in R. *Journal of Statistical Software*, 78(10), 1–39. <https://doi.org/10.18637/jss.v078.i10>
- Madsen, H. (2007). *Time series analysis*. Chapman & Hall. <https://doi.org/10.1201/9781420059687>
- Madsen, H., & Thyregod, P. (2011). *Introduction to general and generalized linear models*. CRC Press.

## Appendix A

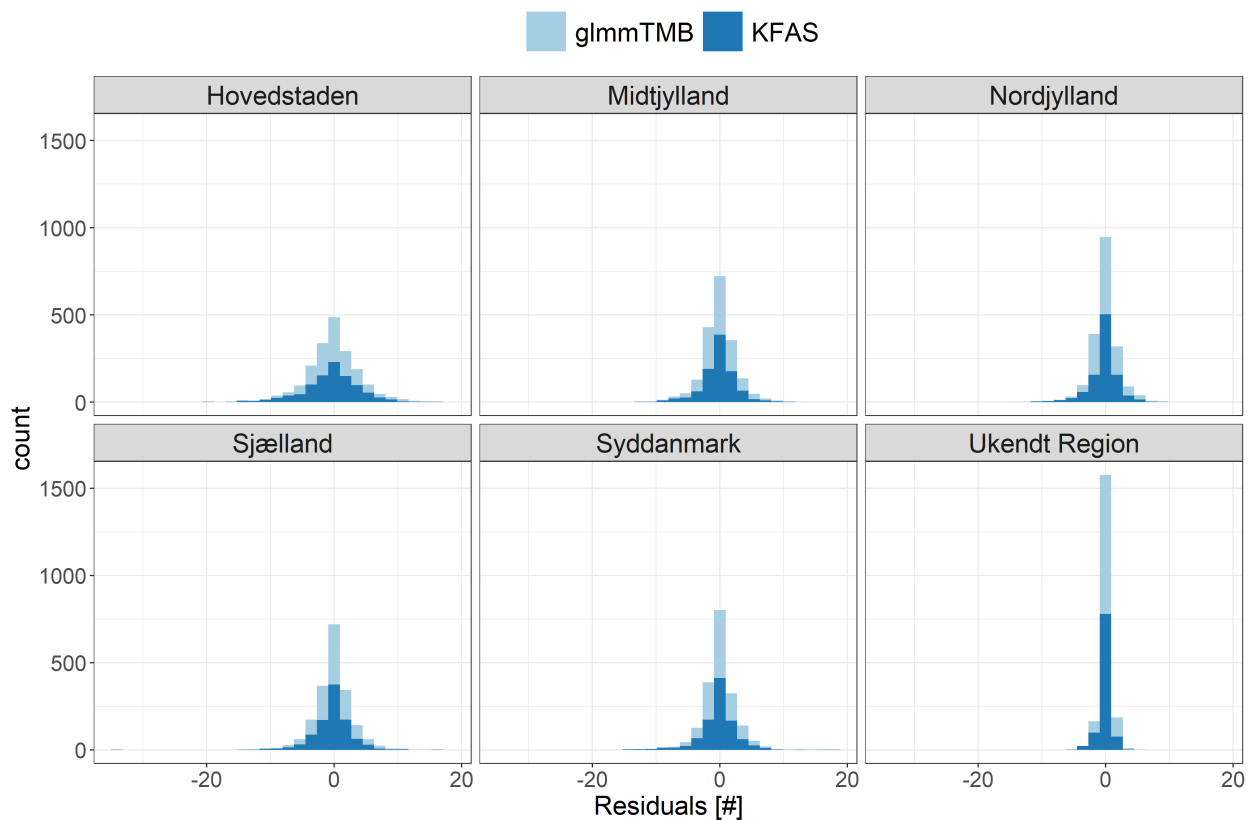


Figure 9: Model residuals from `glmmTMB` and `KFAS`.

## Appendix B

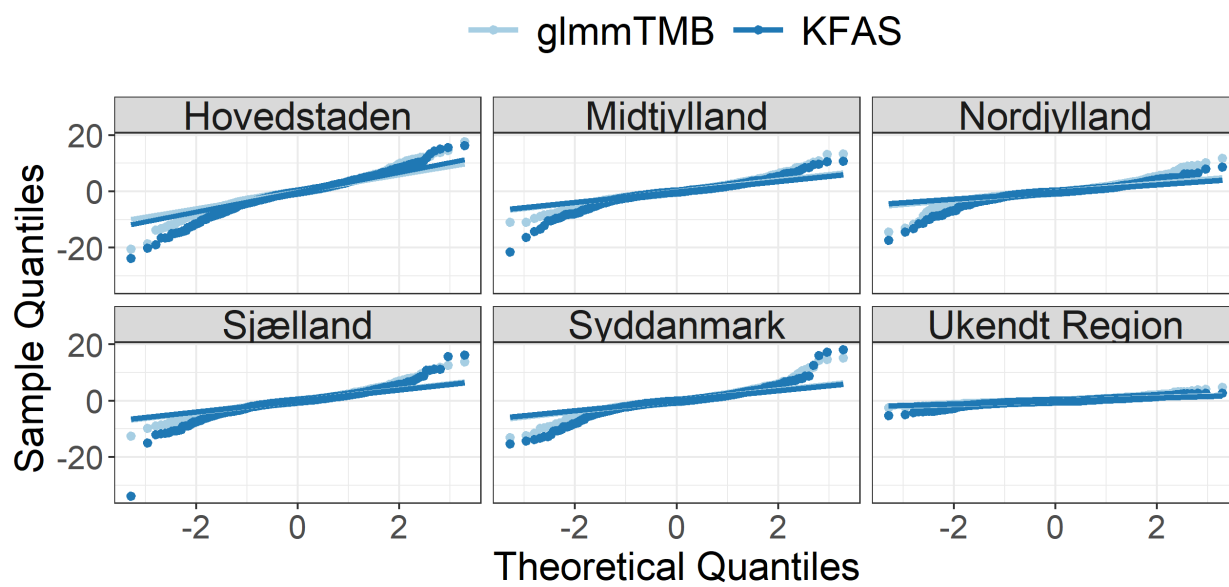


Figure 10: Write me

## Appendix C

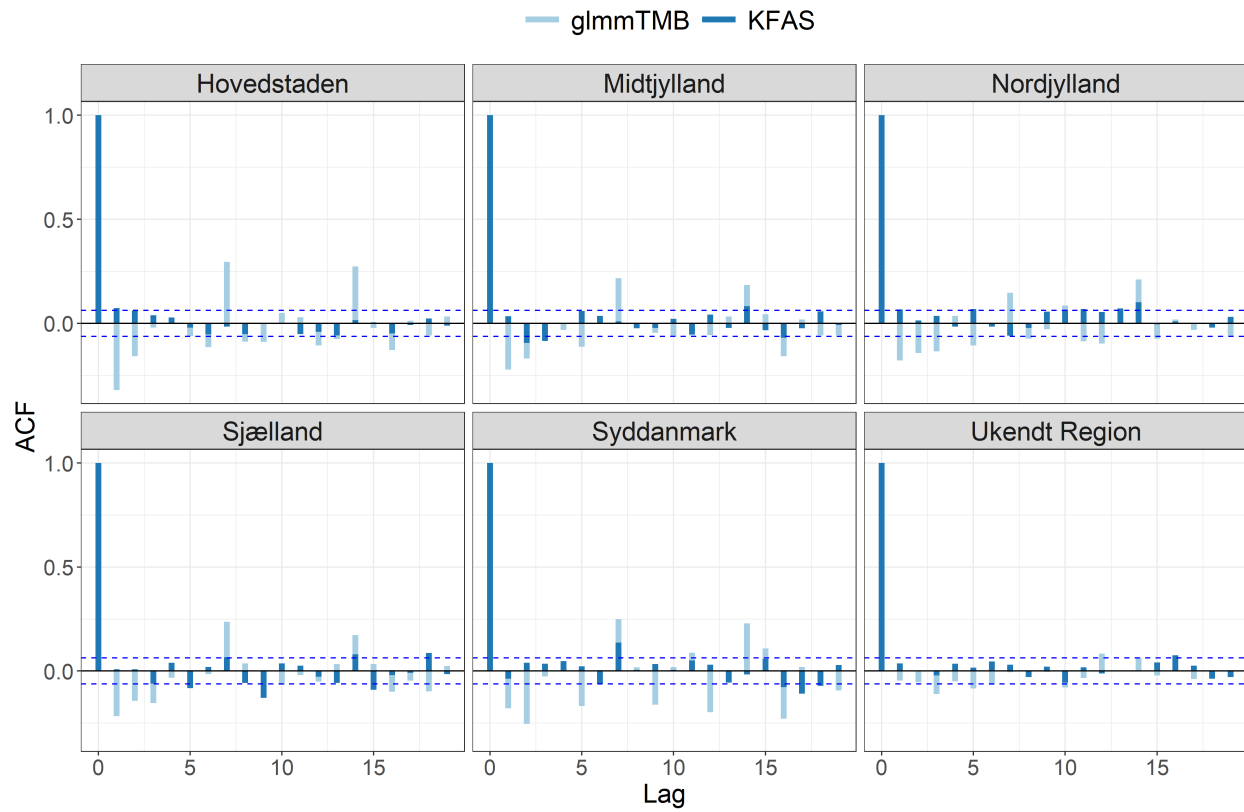


Figure 11: Write me