

Automated and Early Detection of Disease Outbreaks

AEDDO

Master Thesis



Automated and Early Detection of Disease Outbreaks
AEDDO

Master Thesis
August, 2023

By
Kasper Schou Telkamp

Copyright: Reproduction of this publication in whole or in part must include the customary bibliographic citation, including author attribution, report title, etc.

Cover photo: Vibeke Hempler, 2012

Published by: DTU, Department of Applied Mathematics and Computer Science,
Richard Petersens Plads, Building 324, 2800 Kgs. Lyngby Denmark
www.compute.dtu.dk

ISSN: [0000-0000] (electronic version)

ISBN: [000-00-0000-000-0] (electronic version)

ISSN: [0000-0000] (printed version)

ISBN: [000-00-0000-000-0] (printed version)

Approval

This thesis has been prepared over six months at the Section for Dynamical Systems, Department of Applied Mathematics and Computer Science, at the Technical University of Denmark, DTU, in collaboration with Epidemiologisk Forskning / Modelgruppen at Statens Serum Institut, SSI, in partial fulfilment for the degree Master of Science in Engineering, MSc Eng., Quantitative Biology and Disease Modelling.

It is assumed that the reader has a basic knowledge in the areas of statistics.

Kasper Schou Telkamp - s170397

.....
Signature

.....
Date

Abstract

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Acknowledgements

Lasse Engbo Christiansen, Senior Researcher, Statens Serum Institut

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Jan Kloppenborg Møller, Associate Professor, Technical University of Denmark

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Contents

Preface	ii
Abstract	iii
Acknowledgements	iv
1 Introduction	1
2 Literature	3
3 Dataset	5
4 Methods	7
4.1 Hierarchical models	7
5 Results	11
5.1 Case studies	11
5.2 Simulation study	11
6 Discussion	13
7 Conclusion	15
Bibliography	17
A Some probability functions	19
A.1 The Poisson distribution model	19
A.2 The Gamma distribution model	19
A.3 The Negative Binomial distribution model	19
B Proofs	20
B.1 Hierarchical Poisson Gamma model	20

1 Introduction

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

2 Literature

In recent years there has been a surge in interest for statistical methods for automated and early detection of infectious disease outbreaks. The methodologies ranges the spectrum of statistical methods and includes regression techniques, time series methodology, methods inspired by statistical process control, methods incorporating spatial information, and multivariate outbreaks detection. A review of the aforementioned methods can be found in Buckeridge (2007) and Unkel et al. (2012).

Write a paragraph that focus on Farrington et al. (1996) **and** Noufaily et al. (2013)!

Moreover, state-of-the-art methods for aberration detection is presented in Salmon, Schumacher, and Höhle (2016) and implemented in the R package **surveillance**, which is available from the Comprehensive R Archive Network at <https://cran.r-project.org/web/packages/surveillance/index.html>. The R package includes methods such as the Farrington method introduced by Farrington et al. (1996) together with the improvements proposed by Noufaily et al. (2013). As a Bayesian counterpart to these methods the BODA method presented by Manitz and Höhle (2013) allows for easy integration of covariates. Common for these methods are that no accumulation of evidence takes place and they detect aberrations only when the count of the currently monitored timepoint is above the threshold. Another method, which is also implemented in **surveillance** is the negative binomial cumulative sum of Höhle and Paul (2008) that allows for the detection of sustained shifts by accumulating evidence over several timepoints.

3 Dataset

In Denmark the surveillance of infectious diseases is carried out by Statens Serum Institut (SSI). The surveillance is a central part of the national and international disease preparedness. Thanks to The National Board of Health Statutory Order on Physicians' Notification of Infectious Diseases, the quality of the Danish surveillance registers is very high. In the Statutory Order it is stated that a number of diseases¹ are individually notifiable for physicians and general practitioners. The notifications include relevant information on the patient, and are notified on a paper form to the Ministry of Health and to SSI, respectively. A modern surveillance system comprises not only collection and registration of disease data, but also timely and continuous communication of knowledge to authorities responsible for treatment, prevention, and control. The national surveillance system comprises only diseases of serious character, diseases that are particularly infectious, and most of the vaccine-preventable diseases.

For the purpose of this master thesis, only a subset of the individually notifiable diseases are considered. ...**(WRITE SOMETHING ABOUT WHAT CHARACTERISES THE CHOSEN SUBSET I.E. SEASONALITY, INCIDENCE, ETC.)**... and the subset includes: Shiga toxin producing Escherichia coli (STEC), **X, X, X (CHOOSE 4)**.

Include a section with plots and tables of the data.

¹For a full list of individually notifiable diseases see https://www.ssi.dk/sygdomme-beredskab-og-forskning/anmeldelse-af-sygdomme/lovpligtige-meldesystemer/individ_anmeldelses_sygdomme

4 Methods

4.1 Hierarchical models

In this section the hierarchical model and the generalized mixed effect model that is used to model the count observation y_{it} for the specific age groups, $i = 1, \dots, 11$, and in the monthly period ranging from 2008 to 2022, $t = 1, \dots, 180$, is presented. For an introduction to the concept of hierarchical models see Madsen and Thyregod (2011).

4.1.1 Hierarchical Poisson Normal model

The conditional distribution of the count observations, $Y|u$, are assumed to be a Poisson distribution with intensities λ_{it} . Also, we shall assume that the count is proportional to the population size, x_{it} , within each age group, i , at a given time point, t . Hence, in terms of the canonical link for the Poisson distribution the model for the fixed effect is

$$\log(\lambda_{it}) = \mathbf{X}_i^T \beta_{it} + \log(x_{it}) \quad (4.1)$$

Here \mathbf{X}_i is $T \times 6$ -dimensional, and β_{it} contains the corresponding fixed effect parameters. The random effects u_{it} are assumed to be Gaussian

$$u_{it} = \epsilon_{it} \quad (4.2)$$

where $\epsilon_{it} \sim N(0, \sigma^2)$ and are independent and identically distributed, and σ^2 is a model parameter. Henceforth, the model can be formulated as a two-level hierarchical model

$$Y_{it}|u_{it} \sim \text{Pois}(\lambda_{it} \exp(u_{it})) \quad (4.3a)$$

$$u_{it} \sim N(0, \sigma^2) \quad (4.3b)$$

4.1.2 Hierarchical Poisson Gamma model

Likewise, in the compound Poisson-Gamma model the conditional distribution, $Y|u$, of the count observations are assumed to be a Poisson distribution, but this time the intensities, λ_{it} , are defined as

$$\log(\lambda_{it}) = \mathbf{X}_i^T \beta_{it} + \log(x_{it}) \quad (4.4)$$

Here \mathbf{X}_i is $T \times 6$ -dimensional, and β_{it} contains the corresponding fixed effect parameters. Additionally, the random effects u_{it} are assumed to be Gamma distributed. Subsequently, the model can be formulated as a two-level hierarchical model

$$Y_{it}|u_{it} \sim \text{Pois}(\lambda_{it} u_{it}) \quad (4.5a)$$

$$u_{it} \sim G(1/\phi, \phi) \quad (4.5b)$$

In the first stage a random value u_{it} is selected according to a reparameterized Gamma distribution with shape, $1/\phi$, and scale, ϕ . Hence the mean value of the Gamma distribution is 1. Moreover, a fixed effect parameter, λ_{it} , is found for each age group, $i = 1, \dots, 11$.

The Y is generated according to a Poisson distribution with $\lambda_i u_{it}$ as mean value. The the marginal distribution of Y is a negative binomial distribution, $Y \sim \text{NB}(1/\phi, 1/(\lambda\phi + 1))$. The probability function for Y is

$$\begin{aligned}
P[Y = y_i] &= g_Y(y; \lambda, \phi) \\
&= \frac{\lambda^y}{y! \Gamma(1/\phi) \phi^{1/\phi}} \frac{\phi^{y+1/\phi} \Gamma(y + 1/\phi)}{(\lambda\phi + 1)^{y+1/\phi}} \\
&= \frac{\Gamma(y + 1/\phi)}{\Gamma(1/\phi) y!} \frac{1}{(\lambda\phi + 1)^{1/\phi}} \left(\frac{\lambda\phi}{\lambda\phi + 1} \right)^y \\
&= \binom{y + 1/\phi - 1}{y} \frac{1}{(\lambda\phi + 1)^{1/\phi}} \left(\frac{\lambda\phi}{\lambda\phi + 1} \right)^y, \text{ for } y = 0, 1, 2, \dots
\end{aligned} \tag{4.6}$$

where we have used the convention

$$\binom{z}{y} = \frac{\Gamma(z + 1)}{\Gamma(z + 1 - y) y!} \tag{4.7}$$

for z real and y integer values. The marginal distribution of Y is a negative binomial distribution, $Y \sim \text{NB}(1/\phi, 1/(\lambda\phi + 1))$. See proof in B.1.1.

Inference on individual groups

Consider the compound Poisson Gamma model in (4.5), and assume that a value $Y = y$ has been observed. Then the conditional distribution of u for given $Y = y$ is a Gamma distribution

Consider the hierarchical Poisson-Gamma model in (4.5), and assume that a value $Y = y$ has been observed. Then the conditional distribution of u for given $Y = y$ is a Gamma distribution,

$$u|Y = y \sim \text{G}(y + 1/\phi, \phi/(\lambda\phi + 1)) \tag{4.8}$$

with mean

$$\mathbb{E}[u|Y = y] = \frac{y\phi + 1}{\lambda\phi + 1} \tag{4.9}$$

and variance

$$\mathbb{V}[u|Y = y] = \frac{(y\phi^2 + \phi)}{(\lambda\phi + 1)^2} \tag{4.10}$$

Why do we choose the Gamma distribution to represent the variation between days?

The Gamma distribution is chosen for three simple reasons. First of all, the support of the Gamma distribution, $0 < u_{it} < \infty$ conforms to the mean-value space, \mathcal{M} for the Poisson distribution. Secondly, the two-parameter family of Gamma distributions is a rather flexible class of unimodal distribution, ranging from an exponential distribution to a fairly symmetrical distribution on the positive real line. A third reason may be observed in the derivation of the marginal distribution of Y . The fact that the kernel $u^{\alpha-1} \exp(-u/\beta)$ of the

mixing distribution have the same structure as the kernel $u^y \exp(-u)$ of the likelihood function corresponding to the sampling distribution of Y . This feature have the consequence that the integral has a closed form representation in terms of known functions.

4.1.3 Parameter estimation

5 Results

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

5.1 Case studies

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

5.2 Simulation study

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

6 Discussion

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

7 Conclusion

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Bibliography

- Buckeridge, David L. (2007). "Outbreak detection through automated surveillance: A review of the determinants of detection". In: *Journal of Biomedical Informatics* 40.4. Public Health Informatics, pp. 370–379. ISSN: 1532-0464. DOI: <https://doi.org/10.1016/j.jbi.2006.09.003>. URL: <https://www.sciencedirect.com/science/article/pii/S1532046406000980>.
- Unkel, Steffen et al. (2012). "Statistical methods for the prospective detection of infectious disease outbreaks: a review". In: *Journal of the Royal Statistical Society. Series A (Statistics in Society)* 175.1, pp. 49–82. ISSN: 09641998, 1467985X. URL: <http://www.jstor.org/stable/41409708> (visited on 02/15/2023).
- Farrington, C. P. et al. (1996). "A Statistical Algorithm for the Early Detection of Outbreaks of Infectious Disease". In: *Journal of the Royal Statistical Society. Series A (Statistics in Society)* 159.3, pp. 547–563. ISSN: 09641998, 1467985X. URL: <http://www.jstor.org/stable/2983331> (visited on 01/27/2023).
- Noufaily, Angela et al. (Mar. 2013). "An Improved Algorithm for Outbreak Detection in Multiple Surveillance Systems". en. In: *Online Journal of Public Health Informatics* 32.7, pp. 1206–1222.
- Salmon, Maëlle, Dirk Schumacher, and Michael Höhle (2016). "Monitoring Count Time Series in R: Aberration Detection in Public Health Surveillance". In: *Journal of Statistical Software* 70.10, pp. 1–35. DOI: 10.18637/jss.v070.i10. URL: <https://www.jstatsoft.org/index.php/jss/article/view/v070i10>.
- Manitz, Juliane and Michael Höhle (2013). "Bayesian outbreak detection algorithm for monitoring reported cases of campylobacteriosis in Germany". In: *Biometrical Journal* 55.4. Cited by: 18, pp. 509–526. DOI: 10.1002/bimj.201200141. URL: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-84880042358&doi=10.1002%5C%2fbimj.201200141&partnerID=40&md5=4658480d9c1cd84c23335bdb9269f0bc>.
- Höhle, Michael and Michaela Paul (2008). "Count data regression charts for the monitoring of surveillance time series". In: *Computational Statistics and Data Analysis* 52.9. Cited by: 56, pp. 4357–4368. DOI: 10.1016/j.csda.2008.02.015. URL: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-42749087852&doi=10.1016%5C%2fj.csda.2008.02.015&partnerID=40&md5=8d18fdad2a28d99ecbda5441a8fbae3>.
- Madsen, Henrik and Poul Thyregod (2011). *Introduction to general and generalized linear models*. English. Texts in statistical science. CRC Press. ISBN: 9781420091557.

A Some probability functions

This chapter serves as a reference, specifying notation, properties, and moments related to the various distributions used in this master thesis.

Name	Support	Density	$E[Y]$	$V[Y]$
Poisson $\text{Pois}(\lambda)$	$0, 1, 2, \dots$ $\lambda \in \mathbb{R}_+$	$\frac{\lambda^y}{y!} \exp(-\lambda)$	λ	λ
Gamma $G(\alpha, \beta)$	\mathbb{R}_+ $\alpha \in \mathbb{R}_+, \beta \in \mathbb{R}_+$	$\frac{1}{\Gamma(\alpha)\beta} \left(\frac{y}{\beta}\right)^{\alpha-1} \exp(-y/\beta)$	$\alpha\beta$	$\alpha\beta^2$
Neg. Bin. $\text{NB}(r, p)$	$0, 1, 2, \dots$ $r \in \mathbb{R}_+, p \in]0, 1]$	$\binom{r+y-1}{y} p^r (1-p)^y$	$\frac{r(1+p)}{p}$	$\frac{r(1-p)}{p^2}$

Table A.1: Density, support, mean value, and variance for a number of distributions used in this master thesis.

A.1 The Poisson distribution model

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

A.2 The Gamma distribution model

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

A.3 The Negative Binomial distribution model

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

B Proofs

B.1 Hierarchical Poisson Gamma model

This section is a collection of proofs for the derivation of the compound Poisson Gamma model in (4.5).

B.1.1 Probability function for Y

The probability function for the conditional distribution of Y for given u

$$f_{Y|u}(y; \lambda, u) = \frac{(\lambda u)^y}{y!} \exp(-\lambda u) \quad (\text{B.1})$$

and the probability density function for the distribution of u is

$$f_u(u; \phi) = \frac{1}{\phi \Gamma(1/\phi)} \left(\frac{u}{\phi} \right)^{1/\phi-1} \exp(-u/\phi) \quad (\text{B.2})$$

Given (B.1) and (B.2), the probability function for the marginal distribution of Y is determined from

$$\begin{aligned} g_Y(y; \lambda, \phi) &= \int_{u=0}^{\infty} f_{Y|u}(y; \lambda, u) f_u(u; \phi) du \\ &= \int_{u=0}^{\infty} \frac{(\lambda u)^y}{y!} \exp(-\lambda u) \frac{1}{\phi \Gamma(1/\phi)} \left(\frac{u}{\phi} \right)^{1/\phi-1} \exp(-u/\phi) du \\ &= \frac{\lambda^y}{y! \Gamma(1/\phi) \phi^{1/\phi}} \int_{u=0}^{\infty} u^{y+1/\phi-1} \exp(-u(\lambda\phi + 1)/\phi) du \end{aligned} \quad (\text{B.3})$$

In (B.3) it is noted that the integrand is the *kernel* in the probability density function for a Gamma distribution, $G(y + 1/\phi, \phi/(\lambda\phi + 1))$. As the integral of the density shall equal one, we find by adjusting the norming constant that

$$\int_{u=0}^{\infty} u^{y+1/\phi-1} \exp\left(-u(\phi/(\lambda\phi + 1))\right) du = \frac{\phi^{y+1/\phi} \Gamma(y + 1/\phi)}{(\lambda\phi + 1)^{y+1/\phi}} \quad (\text{B.4})$$

and then (4.6) follows

B.1.2 Conditional distribution of Y

The conditional distribution is found using Bayes Theorem

$$\begin{aligned} g_u(u|Y = y) &= \frac{f_{y,u}(y, u)}{g_Y(y; \lambda, \phi)} \\ &= \frac{f_{y|u}(y; u) g_u(u)}{g_Y(y; \lambda, \phi)} \\ &= \frac{1}{g_Y(y; \lambda, \phi)} \left(\frac{(\lambda u)^y}{y!} \exp(-\lambda u) \frac{1}{\phi \Gamma(1/\phi)} \left(\frac{u}{\phi} \right)^{1/\phi-1} \exp(-u/\phi) \right) \\ &\propto u^{y+1/\phi-1} \exp(-u(\lambda\phi + 1)/\phi) \end{aligned} \quad (\text{B.5})$$

We identify the *kernel* of the probability density function

$$u^{y+1/\phi-1} \exp(-u(\lambda\phi + 1)/\phi) \quad (\text{B.6})$$

as the kernel of a Gamma distribution, $G(y + 1/\phi, \phi/(\lambda\phi + 1))$

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Technical
University of
Denmark

Richard Petersens Plads, Building 324
2800 Kgs. Lyngby
Tlf. 4525 1700

www.compute.dtu.dk