

# CANSSI - Letter of Submission

## Summary of changes

The Aims have been revised to emphasize the methodological aspects of the research, with the new aims being as follows:

1. develop bcGAIM, a Bayesian inference methodology for high dimensional cGAIM's in case/crossover models;
2. create an efficient, non-iterative computational algorithm for bcGAIM's based on Laplace approximations;
3. develop non-parametric forms of the dose-response effect which encourage or enforce monotonicity; and
4. engage in interdisciplinary and applied research projects with our subject-area collaborators.

A principle challenge will be to parametrize the model in such a way as to clearly separate the most strongly and most weakly identifiable components, thereby enabling an MCMC implementation which is able to explore the 'weakly identified' parts of the parameter space while also showing where subjective expert-based prior information is needed. Computational challenges will also exist, and the key will be efficiently parallelizing the algorithm (possibly across multiple computers). While the applications of these methods are important and motivate the research, they are confined to a single aim and much of this work will be driven by the group's collaborators in the health sciences.

## SAC Review

1. *Provide details about the mentoring of the HQP, linking HQP directly to specific mentors. Include details on the mentoring of the student at Laval and how the student will be a full participant in the team. Keep in mind that it is unusual for a postdoctoral fellow to play a major role in mentoring, without the involvement of a senior researcher.*

The HQP budgeted in the project are 1 postdoc fellow (PDF), 1 PhD student, 1 MSc student and 2 undergraduate trainees per year. The PDF, likely Kamal Rai, will be based initially in Toronto, supervised by Patrick Brown and Meredith Franklin, and later in Halifax under the supervision of Cindy Feng and Daniel Rainham. The PhD student can be based at any of the three institutions but preference will be given to a candidate at INRS in Quebec. The main advisor of the MSc student will be Cindy Feng and will be located in Halifax. PDF and PhD student will provide day-to-day supervision of the undergraduate trainees, under the guidance of an Investigator.

2. *Clarify how the team will function and how team members will interact and collaborate. Describe the roles of the lead investigators in the project and research.*

The PIs have the main role of the HQP supervision. Each PI acts as a supervisor of one main HQP (PDF, PhD and MSc) and one undergraduate trainee, and also as a co-supervisor for a second HQP. Even though their contribution is essential for the project, the collaborators will play a complementary role in the supervision of the HQP.

Regular meeting (weekly) between each HQP and the associated main advisor, and monthly (or as needed) with the co-advisors. Each month, all team members (PIs, HQP, collaborators) involved in a given sub-project will meet (in person or virtual) to discuss that specific sub-project. Bimonthly all

team members will meet for updates, discussing issues, presenting results and methods, depending on the stage of the whole project (e.g. at the beginning discussions more about hiring, getting data, but at the end more about results).

Trainees will be encouraged to spend at least one month at one of the other two institutions at some point. All members will meet annually, likely at the SSC annual conference.

3. *Provide more details on how the aims associated with the application are related to the statistical methodology.*

The Aims have been revised to focus on methodological contributions, and how the applications motivate the statistical work is more clearly explained. The practical problem to be addressed is assessing how multiple air quality and environmental risk factors (at different time lags) affect short term mortality when the effects are non-linear and the variables strongly correlated. The bcGAIM methodology proposed will handle multiple highly correlation air pollutants and relate them to short term morbidity/mortality with case-crossover models.

4. *Briefly describe the available data sets and the confounders present therein.*

A number of datasets from Canada, India, the US and China are currently being used by the four PI's. Each problem will involve environmental monitoring data of various sorts, and health outcomes consisting of daily case counts. The case/crossover model is in effect adjusting for all individual-level covariates using a partial likelihood, and only cause-specific georeferenced case counts are required.

Canadian daily air pollution data has been acquired from the National Air Pollution Surveillance (NAPS) Program, a network of 250 stations across Canada that is managed by Environment Canada and Environment and Climate Change Canada. Data on important environmental confounders of air pollution, such as temperature, wind speed and humidity will be acquired from weather stations operated by Environment and Climate Change Canada.

Canadian daily health outcome data, which are not publicly available, have been provided to the investigators by Health Canada and INSPQ (the Quebec Public Health Institute). These data provide daily counts of cause-specific hospital admissions and mortality at the census tract or city level. Confounders in the health dataset include age, gender, race, and various markers of socio-economic status at the tract/city level.

COVID-19 data, which are publicly available, will be acquired from provincial authorities. We will focus on building models for COVID-19 mortality as they are more reliable than the case data, which may suffer from bias since not everyone was tested due to limited testing capacity at the beginning of the pandemic.

The multi-pollutant cause-specific morbidity and mortality models will be developed separately from the multi-pollutant COVID-19 models. Both the nature of the data (count vs individual) and the available confounders are different between the two health datasets requiring different applications of the bcGAIM approach.

5. *Discuss why you model mortality from asthma, rather than occurrence.*

We are currently looking at both mortality and hospitalizations for a number of conditions, and bcGAIM can be fit to any health outcome of interest. The outcomes will be chosen in consultation with epidemiologists with whom we are collaborating. However, we do have a few outcomes in mind and they will be chosen based on the application of bcGAIM. For example, with the Health Canada data we will look at cause-specific morbidity and mortality including cardiopulmonary outcomes that have been examined in our previous studies (Franklin et al 2007,2008; Zanobetti et al 2009). To clarify the SAC's question, asthma will only be examined as a morbidity (i.e. occurrence), not as a cause of mortality.

As stated in the previous comment, COVID-19 mortality will be of interest due to issues with under-reporting of cases when tests were not available.

6. *Detail the technical challenges to be addressed, provide references to what is already known and clearly state what is to be developed. In your discussion, include the following.*

- (a) *Provide explicit information on the kinds of constraints you will place on the prior distributions. Provide more extensive references to the literature on shape constrained models.*

We now elaborate on the shape constrained non-parameteric effects. We will consider three approaches: basis function methods along the lines of @ramsay1998, reparameterizing Gaussian processes to accomodate shape-constrained priors, and various forms of random walks (i.e. higher order, or RW1 + drift).

- (b) *Describe the challenges of introducing Bayesian approaches.*

There are three main challenges to be addressed: finding a parametrization of bcGAIM which clearly separates combinations of parameters which are well identified and those which are weakly identified by the data; creating an MCMC algorithm which mixes well when the dimensionality is large and the likelihood surface flat; developing an INLA-like algorithm which is more easily automated than MCMC. We now elaborate on these challenges in the proposal.

- (c) *You propose to initiate the models in STAN. Does this mean the research is very straightforward? If not, what are the challenges?*

It is unlikely that the most straightforward parametrization of bcGAIM will work well in Stan when the dimensionality is high. Finding an efficient implementation in Stan and developing alternatives to Stan are two of the key challenges. Developing a non-parametric model which is monotone (or encourages monotonicity) is another challenge.

## Reviewer 1

1. *From a methodological perspective, the proposed extensions are not particularly novel. However, the application of Bayesian nonparametric regression in the context of air pollution epidemiology is novel.*

The *Research Aims* contains additional details on the proposed extensions in the bcGAIM. Of these, the "nested model" construction of shape-constrained priors for Gaussian processes and the approximations proposed to render the hierarchical model less computationally demanding are perhaps more novel than they may have appeared in the LOI.

## Reviewer 2

1. *What sort of research would be needed to construct the new bcAQHI once the new bcGAIM is built and computationally implemented?*

The bcGAIM will output a relative risk for every combination of (measured or forecasted) pollutant values input into the model. The bcAQHI convert these relative risk distributions into warnings, based on posterior probabilities of being above cut points or risk thresholds. The key task is determining biologically relevant cutpoints, which we may do in a model-based way or using expert knowledge. Finally, forecasted pollutant values for next-day predictions will be provided by Environment Canada.

2. *The cGAIM itself is by no means a new idea. The original model for a single group ( $K=1$ ) goes back to Hardle (1993). In fact, one seems to get Hardle's model if one drops the  $f_k$ 's from the model. Wang et. al. (2015) presents a multigroup version ( $K>1$ ) to get around the curse of dimensionality, the whole point of this approach. But special cases were published between 1993 and 2015.*

There are new features in the cGAIM – it considers constraints and groupwise additive index terms, while much of the existing literature only considers one or the other. While Hardle, Hall, and Ichimura (1993) examine a single index model and Wang et. al. (2015) consider a multiple index model, neither consider constraints. Two papers that consider constrained estimation are Xia and Tong (2006), where the authors constrain  $s$  to be monotonic and the components of  $\alpha$  to be non-decreasing, and Fawzi et al. (2016), where the authors constrain the components of  $\alpha$  to be non-negative and sum to one but do not constrain  $s$ . In comparison, the cGAIM allows any linear constraint to be placed on  $\alpha$  and different shape constraints on  $s$  including monotonicity, convexity, and concavity (Masselot et al., 2020). Additional comparisons between the cGAIM and bcGAIM are given in the *Research Aims* section.

3. *The third main topic seems the most novel in as much as it will show how the new bcAQHI might be used to assess COVID-19 mortality. Of course, it would seem more reasonable to me to build a new bcAQHI designed specifically for that purpose. And that led me to wonder about the health outcome to be used to fit the bcGAIM-the all-causes mortality-to get the alpha and in turn the index. For example, ozone would seem more relevant as a risk index for asthma and PM 2.5 for COPD. Why not just publish the pollutant concentrations themselves?*

It is the bcGAIM, rather than the bcAQHI, which will be used with COVID-19 mortality. A covid-specific bcAQHI will be computed, the revised application is now clearer on this.

Regarding reporting pollutant concentrations, in Canada they are publicly available via the National Air Pollution Surveillance (NAPS) Program. However, it is very difficult to understand their health effects without using a model-based approach. For the multi-pollutant model, the key benefit of the bcGAIM (and the bcAQHI) is that it provides an ease of interpretation for these health effects that just publishing data does not. Compared to single-pollutant models, it provides a measure of the relative risk of the mixture of air pollutants in the ambient air. Moreover, this is a much better representation of the health risks an individual will face than what is estimated by a single pollutant model.

4. *The conversion of the cGIAM to a bcGIAM should be feasible. A major challenge will be the big data problem since daily health counts for all-cause mortality will be modeled, and this using a semi-parametric model. But the task would be simplified by the intended approach of designing a different model for each Canadian city, how many we don't know. But they do not intend to incorporate random city effects in the Bayesian framework to enable strength to be borrowed deficiency in the intended approach, but perhaps a compromise needed for feasibility.*

There are significant challenges relating to the size of the data. We will be examining 25+ regions in Canada over a 20+ year period (6,000+ days) of daily data. Fitting the bcGAIM to each city independently is computationally feasible, while fitting a hierarchical model across 25 regions is more challenging but perhaps more desirable. We have given further thought to this goal since submitting the LOI, and will be implementing various Laplace approximations in the target density to lessen the computational burden. This would allow us to produce estimates of the relative risks of the pollutant

mixture using a hierarchical model. The details of the proposed approximation are discussed in the *Research Aims* section.

5. *A complicated issue and I had to do a lot of digging to figure out how the team of Investigators was assembled. In part, this involves those listed as Collaborators. The LOI could have done a better job of clarifying the links of the Investigators and Collaborators in the proposed project.*

This has been cleared up. In Quebec city, Fateh Chebana is the PI and Céline Campagna is collaborator from INSPQ. In Halifax, Cindy Feng is the PI and Daniel Rainham is collaborator. In Toronto the PIs Patrick Brown and Meredith Franklin are at the University of Toronto, the former is cross-appointed at Centre for Global Health Research, St. Michael's Hospital. Prabhat Jha at CGHR and Hwashin Shin at Health Canada are the two collaborators connected to the Toronto group.

6. *It is challenging to coordinate and run such a program successfully and we don't get a clear impression from the LOI that the applicants have thought about this issue very much. What is clear is that it is designed to provide the pipeline from data through to the AQHI. That is excellent. What is not clear is how the collaboration is to be managed and what I for one would like to see if a proposal is invited is an active group interaction plan.*

Indeed, such of multidisciplinary and multi-location projects are challenging but interesting. This is particularly true in our project where members are at least in three provinces (Ontario, Quebec and Nova-Scotia) and from different disciplines (including statistics, epidemiology, environment, public health). Some PIs have already been involved in similar projects (e.g. Cindy Feng in a previous CRT by CANSSI) and as leader (e.g. Fateh Chebana with a major project with INSPQ). Hence, PIs have experience in successfully managing such kind of projects. In addition, in this specific project, the clarity of the role distributions, the different levels of contributions and responsibilities, the planned meetings (at different locations, for different purposes, etc), the complementarity between all team members, all are ingredient to project successful.

7. *The discussion of the bcGIAM is well done including the parts about the implementation of the computation strategies. But the LOI should have explained it better. For a start the big picture with its three main objectives should have been better described in general terms and how they relate to one another before diving into details re the cGIAM.*

The aims of the project and the research plan have been more clearly laid out, we hope the reviewers will find this revision much clearer.

## Reviewer 3

1. *The naive case fatality rate that the investigators plan to explore with regard to Objective 3 is subject to errors caused by an undercount of both the numerator and the denominator and known to be a poor measure of the mortality risk of the disease.*

Correct, we will consider two forms of denominators when modelling COVID-19 mortality risk. The population fatality rate will be used initially, using census-based age-sex stratified population counts. Improved estimates of the number of covid cases should soon be available through studies involving antibody data, including the ABC study which Dr. Brown is involved in. Carried out in conjunction with the Angus Reid Forum, this study is repeatedly testing individuals for antibodies and the number of undetected cases will be estimated.

2. *Although four methodological advancements of the bcGAIM are listed in the Methods section, those seem to be rather incremental changes of the cGAIM.*

The *Research Aims* section has additional material explaining the differences between the cGAIM and bcGAIM. The major benefit is in quantifying the uncertainty of  $\alpha$ . Secondary benefits include being able to specify the strength of the monotonicity constraint, and the ability to extend the bcGAIM to additional smooth functions  $s$  or additional covariates in each  $s$ . As detailed in *Research Aims*, quantifying the uncertainty in the estimate of  $\alpha$  is crucial for conducting inference with the multi-pollutant model. The extensibility of the bcGAIM and being able to control the strength of the shape constraint are also important features that will play an important role when in developing the multi-pollutant model.

3. *How the bcGAIM helps with developing a simple, intuitive air quality index that simultaneously accounts for the health effects of multiple air pollutants is not explained.*

This was clarified in the response to the first question by Reviewer 2.

4. *The potential for impact in statistics and inferential data science seems to be marginal considering that methodological advancements of the bcGAIM listed in the Methods section are deemed to be rather incremental changes of the cGAIM.*

We respectfully disagree. The bcGAIM has novel innovations that will enable new statistical reasoning to be applied to mixtures of interest, such as in the multi-pollutant problem. There are a number of statistical challenges in developing the bcGAIM. For one, fitting a model to daily observations across 25+ regions in Canada is a significant computational task. The bcGAIM attempts to address this by applying a (non-linear) one-dimensional function to a linear combination of related covariates, which both eases the computational burden and improves the interpretability of the model. To accomplish this, we must set shape-constrained priors on  $s$  and estimate  $\alpha$ . These are both challenging problems and are discussed in more detail in the *Research Aims* section. It also has additional material to clarify the differences between the cGAIM and bcGAIM.

## Reviewer 4

1. *The team has the potential to provide an excellent environment for interdisciplinary training of students. It would be helpful to add the names of the faculty supervisors/collaborators to the mentoring plan. The part for the roles of personnel could be further clarified.*

Clarification can be found in the answer to the second question posed by Reviewer 2, as well as the *Anticipated Roles of Trainees* and *Anticipated Organization of Collaboration* sections. These two sections contain more details on how collaboration is organized, and the role of the investigators.



## Reviewer 5

1. *A potential weakness would be that while well motivated by multi-pollutant modeling, the investigators could strengthen the proposal by identifying other applications where this form of model would be applicable.*

The *Research Aims* section contains a brief discussion of other applications, with some further discussion in *Anticipated Roles of Trainees*.

2. *What is the relationship between the linear combination that goes into the smooth function and the air quality index? How will the index provide measures that indicate it is safe or not based on the smooth function and or linear combination of the exposures? In terms of data, is this based on hospital admissions (asthma or other conditions or only mortality). More details on data sources related to aims would be helpful.*

The bcGAIM estimates the relative risk of the linear combination of the pollutants, and is translated into an air quality index based on cutpoints on the levels of the relative risk. For the air quality index, we are collaborating with epidemiologists and other scientists and intend to look at many different health outcomes. We have daily cause-specific case counts for 25 cities in Canada, and very detailed mortality data from the Million Deaths Study carried out in India.

3. *While the researchers describe this as a constrained or shape constrained model, the proposal lacks details about what shape constraints are desired. There is a significant literature on Bayesian shape constrained modeling (monotonicity in particular) but limited references are provided. How does the proposed research build on this and what will be novel and contribute to general statistical methodology? Is this a straightforward model to fit in STAN or code directly or are there methodological advancements to be made there?*

The *Research Aims* section contains additional details on the bcGAIM, more extensive references to the literature on shape-constrained inference, and details on the challenges to fitting the bcGAIM in Stan. We do not expect Stan to work well without significant reparametrizations and optimizations, and developing an alternative to Stan is one of the aims of the project.

4. *The researchers cite Stringer et. al. (2020) as developing a Bayesian single pollutant version of a case-crossover model using non-MCMC methods such as Integrated Nested Laplace Approximations (INLA). Given the space limitations of the proposal it is not clear that such results will immediately carry over to the Bayesian shape constrained model where both the (constrained) smooth function of the linear combination of the exposures, the weights in the linear combination of exposures and the smooth functions of confounders has to be estimated, in addition to the other smooth functions of confounders. INLA like methods have been used to solve an array of complex problems, so this may be feasible. While this might provide computational efficiencies in point estimation, how does this address the importance of uncertainty quantification of the Bayesian model over the frequentist model of Masselot et. al. (2020)?*

The *Research Aims* section gives more detail on the INLA-like approximation. There will certainly be challenges in using the methods from Stringer et al 2000 with bcGAIM due to the much higher dimension of the non-linear parameters. Given any fixed values of the weights  $\alpha$ , the INLA-like methodology from Stringer will compute posteriors. Allowing for uncertainty in  $\alpha$  will require running this algorithm many times for different values of  $\alpha$ , finding an efficient way to do this is one challenge. Adapting the method to handle the constrained form of the random effect is another important issue that must be addressed.

5. *In discussing priors to induce shape constraints, the proposal rejects the idea of placing a prior on the expansion of the smooth function, but rather to place it directly on  $s$ . What types of priors on functions spaces are going to be used; Gaussian Process priors or others? How are the constraints incorporated?*

The *Research Aims* section contains a brief discussion of shape-constrained Bayesian inference, and additional discussion on our approach to developing shape-constrained priors. Several different methods will be considered.

6. *While the impact of pollutant exposure and COVID-19 exposure is suggested by Wu et. al. (2020), how does this model relate to the models for total mortality? Will this be a joint model for COVID versus non COVID deaths with common smooth s function or different smoothed functions? Or will there be different linear combinations of exposures or the same? Or will this utilize the proposed Air-Quality index? What data are available and do they provide the necessary information about potential confounders or other covariates? (socio-economic status, access to health insurance, housing status (group living such as nursing homes, dorms, single family, number of family members etc) co-morbidities, etc. Are individual level data available or is this aggregated data at say a county level? Missing data is clearly an issue with COVID-19 deaths but is not addressed.*

We will fit the bcGAM model to COVID-19 deaths separately from other health outcomes, finding a linear combination specific to COVID. We are focusing on temporal variation in deaths and pollution, the case/crossover model is in effect controlling for any covariates which are not time varying (i.e. socio-economic variables). COVID-19 death counts are available at a the public health unit level in Canada, county level in the US, and reasonably small areas for most countries in Europe. Since submitting the LOI, more detailed COVID-19 death data has been made available. For example the Government of Ontario's Treasury Board Secretariat provides many COVID-19 data sets at the individual level, including "Long-Term Care Home COVID-19 Data", "Confirmed positive cases of COVID-19 in Ontario", and "Status of COVID-19 cases in Ontario". The "Confirmed positive cases of COVID-19 in Ontario" data set contains age, gender, location (by public health unit), and the patient outcome. The "Status of COVID-19 cases in Ontario" data set contains daily tests completed, test outcomes, case outcomes, current hospitalizations, and current patients in ICUs. Other provinces only provide data at the census tract level, with breakdowns for covariates such as age and gender.

It is important to note that none of these data sets were available at the time of our LOI submission. The Government of Ontario has gradually made additional information available to the public, a trend that will likely continue in Ontario and other reporting regions. It's very reasonable to expect more data in Ontario (and elsewhere) to be made available over time.

7. *As this is a very application motivated proposal it would be useful to know what data are available for each of the aims and how they will be integrated.*

The applicants and their collaborators have extensive experience working with data from a variety of sources. See response 4 to the SAC comments.

8. *The proposal mentions developing random effects in the smooth function. What do these capture and what is the motivation? i.e spatial random effects, treating the weights as random effects (to allow spatial variation). Additional clarity would be helpful.*

By random effects, we are referring to Gaussian processes such as random walks. These are alternatives to spline functions for non-parametric models, the language around random effects has been cleaned up.