

Bayesian Regression Methods

Motivation

We are interested in using Bayesian regression techniques to characterize the nature of complex interactions in exposure mixture studies. We begin by reviewing definitions for what constitutes a complex interaction and why interactions are relevant from public health and biological relevance.

Interactions from a statistical perspective

First, we define additivity and non-additivity in the traditional statistical paradigm [siemiatiycki_biological_1981]. Suppose we have two variables x_1 and x_2 , and we want to consider their effect on some outcome of interest. If specifying [effect due to x_1 and x_2] = [effect due to x_1] + [effect due to x_2] can adequately capture this relationship, then we say that x_1 and x_2 each have an **additive effect** on the outcome and that there is no interaction between them. On the other hand, if there is variability in the outcome that can be captured by an additional term equal to some function of x_1 and x_2 , we say that there is a **non-additive interaction** between x_1 and x_2 . In this case, [effect due to x_1 and x_2] = [effect due to x_1] + [effect due to x_2] + [effect due to $f(x_1, x_2)$], where f is non-zero for some values of x_1 and x_2 .

For our sake, we consider any non-additive interaction to be complex, meaning that they are difficult to detect. To see why, let us consider running a linear regression for Y on x_1 and x_2 . The regression would be defined as

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \hat{\beta}_{12} f(x_1, x_2),$$

where the β 's represent the effect sizes. We can see that the form of the interaction must be explicitly specified in the formulation of the model. Most commonly, a multiplicative interaction is assessed, where $f(x_1, x_2) = x_1 * x_2$. However, a non-additive interaction can take on many different forms, the true nature of which is difficult to determine analytically.

We used a two-predictor case above, but interactions can also exist between two or more variables (i.e., two-way by $f(x_1, x_2)$, three-way by $f(x_1, x_2, x_3)$, etc.). And if we wanted to assess all possible interactions, the number to consider quickly becomes intractable in high-dimensional settings. For instance, consider modelling 10 predictors in the above linear regression setting. In order to be assessed, each interaction must be explicitly specified as a new term in the model. Even if we only considered one form for the interaction, including all possible two-way interactions would involve adding $\binom{10}{2} = 45$ additional terms to the model, and all possible three-way interactions would add $\binom{10}{3} = 120$ additional terms.

It is important also to acknowledge, here, that there is a limit to how many variables can be included in an interaction before it becomes incomprehensible to most humans. For instance, halford_how_2005 suggest that there is a steep decline in interpretability from three- to four-way interactions, and that five-way interactions are only interpreted correctly at chance level [halford_how_2005]. Hence, for practical purposes, we will limit our exploration to two- and three-way interactions.

Mechanistic and public health relevance

Thus far, we have discussed interactions within a statistical paradigm. However, in addition to being an interesting estimation challenge, non-additive interactions are also relevant in exposure mixture studies from both a mechanistic and public health point of view.

From a mechanistic perspective, a non-additive statistical interaction between two chemical exposures suggests that these molecules may be functionally interacting with each other. Theoretical models propose that such interactions can be classified as either synergistic or antagonistic [heys_risk_2016; plackett_quantal_1952]. In a synergistic interaction, the joint effects of a mixture exceed the independent effects

of each component. This usually occurs if a chemical induces an enzyme involved with the activation of a second chemical or if a chemical inhibits an enzyme that would have otherwise degraded a second chemical. For example, various synergistic interactions between heavy metals have been documented in the literature [e.g., @paithankar_heavy_2021; @wang_associations_2018].

On the other hand, in an antagonistic interaction, the joint effects of a mixture are less than their independent effects. This can occur either through competition at the target site of an enzyme or through direct chemical reactions with each other. In general, synergistic interactions are more concerning in risk assessments, as they lead to underestimation of the true toxicity of a mixture.

It should be noted, though, that while statistical interactions may provide some insight into how exposure mixtures are related to health, they cannot confirm their underlying biology [@vanderweele_tutorial_2014]. If the goal is to assess a meaningful biological interaction, then the discovery of a statistical interaction should be followed up by a functional study.

Now, from a public health perspective, we might be interested in how exposure mixtures interact with other covariates, or, in other words, how social and health factors might mediate the relationship between a health outcome and chemical exposures [@vanderweele_tutorial_2014]. In our case, we can include these additional covariates in the exposure mixture model, where, statistically, they would contribute to the model in the same manner as another chemical exposure: a predictor. A statistical interaction in our model between a covariate and an exposure would indicate that the *magnitude* of the effect of reducing the level of an exposure might differ across various levels of the covariate. This finding could be relevant to public health policy makers, as the potential benefit of regulating a pollutant might differ across groups. For instance, it has been suggested that nutritional intake may modify susceptibility to chemical exposures [e.g., @kannan_exposures_2006; @kordas_interactions_2007].

In many cases, we might assess a covariate related to health inequity, such as socioeconomic status. We provide a cautionary comment, here, that an interaction term should not be the *sole* measure used to measure a health disparity [@ward_how_2019]. In this case, we should first consider the independent, additive association between the covariate and levels of exposure or rates of a health outcome, in order to contextualize the meaning of a potential interaction term.

Bayesian kernel machine regression

In this section, we introduce the theory of BKMR. First, we define the notation that we will be using for kernel machine regression:

- X_m is a predictor variable in the predictor matrix \mathbf{X} with $m = 1, \dots, M$, measuring exposure variables or covariates
- \mathbf{x}_i is a vector of values for a single observation in \mathbf{X} with $i = 1, \dots, n$
- x_{im} is the i th observation of X_m
- Y_i is an observation of \mathbf{Y} , measuring the health outcome in this case
- $h(\cdot)$ is the flexible function relating \mathbf{x} to \mathbf{Y}
- k is the kernel function, the Gaussian in this case
- \mathbf{K} is the $n \times n$ kernel matrix, with (i, j) th element $k(\mathbf{x}_i, \mathbf{x}_j)$
- ρ is the parameter inside the kernel function which controls smoothness
- τ is the parameter multiplied by the kernel matrix to relate \mathbf{K} to h , and
- $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$ are the residuals of the response.

And, we define the notation that we will be using specific to BKMR:

- $r_m = 1/\rho_m$ is an augmented variable in \mathbf{r} in the kernel matrix, controlling smoothness
- δ_m is an indicator variable in $\boldsymbol{\delta}$ which represents inclusion in the model
- \mathcal{S}_g is a group of partitioned predictors with $g = 1, \dots, G$

- $\{\delta_m | \mathbf{x}_m \in \mathcal{S}_g\}$ is an indicator variable in $\delta_{\mathcal{S}_g}$ which represents inclusion of a parameter in group g in the model
- π is the prior probability of inclusion of a predictor in the model, and
- $\lambda \equiv \tau\sigma^{-2}$ is a convenient way to define the prior on τ .

Kernel machine regression

We begin by introducing kernel machine regression, with attention to its specific implementation in BKMR. First proposed by @nadaraya_estimating_1964 and @watson_smooth_1964, kernel machine regression is a nonparametric regression technique that can be used to capture non-linear effects and non-additive interactions. In this introduction, we follow the presentation of kernel machine regression provided by @bobb_bayesian_2015.

To contextualize this method, we start at the typical linear regression setting,

$$Y_i = \mathbf{x}_i \boldsymbol{\beta} + \epsilon_i,$$

where Y_i measures a health outcome at a given point, $\mathbf{x}_i = [x_{i1}, \dots, x_{iM}]$ is a vector of M exposures or covariates (hereafter referred to as predictors), $\boldsymbol{\beta}$ is a vector of weights, and ϵ_i is a random variable from $\epsilon \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$. We can see that this function assumes that there is a linear relationship between the exposure and the response, and that the combined effects of multiple exposures are additive.

Kernel machine regression defines this relationship using a flexible function $h : \mathbb{R}^M \rightarrow \mathbb{R}$, where

$$Y_i = h(\mathbf{x}_i) + \epsilon_i.$$

Here, $h(\cdot)$ is represented by the function $k(\cdot, \cdot)$, a kernel. The kernel controls the covariance, or the similarity, between values of $h(\mathbf{x})$ and as such ensures that points near each other on the prediction surface will have similar values — or, in other words, that the prediction surface will be smooth. In the case of kernel machine regression, we define a positive definite kernel where $k : \mathbb{R}^M \times \mathbb{R}^M \rightarrow \mathbb{R}$.

There are many choices of functions for k . BKMR uses the Gaussian kernel, also known as the radial basis function or, sometimes, the squared exponential kernel. The Gaussian kernel is defined as

$$k(\mathbf{x}, \mathbf{x}') = \exp \left\{ - \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\rho} \right\},$$

where $\|\mathbf{x} - \mathbf{x}'\|^2 = \sum_{m=1}^M (x_m - x'_m)^2$ for a set of predictors values \mathbf{x} and the predictor values of another subject \mathbf{x}' , and ρ is a tuning parameter that controls the relationship between the correlation between two points and their distance. Greater values of ρ will enforce more dependence between points and make the resulting function smoother. h is related to k by a multiplicative constant τ , a tuning parameter which controls the vertical scale of h .

Now that we have defined h and k , we can think about how to characterize the relationship between our response and predictors. Kernel machine regression is a nonparametric technique because it does not specify a functional form for this relationship. Hence, we will think about estimating the response at a particular query point. Operationally, @muller_weighted_1987 demonstrated that kernel machine regression uses a weighted average of all the observations in the dataset to estimate the response, defined as

$$\bar{Y} = \frac{\sum_{i=1}^n w_i Y_i}{\sum_{i=1}^n w_i},$$

with some set of weights $\{w_i\}_{i=1}^n$. Intuitively, we want to weight the observations that are closer to the query point more heavily. Using the Gaussian kernel as a weight allows us to achieve this. Replacing the weight with the Gaussian kernel, we get

$$\bar{Y} = \frac{\sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i) Y_i}{\sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i)}.$$

As we move through the predictor space, we can think of the prediction as a continuous moving average of local points in the dataset. The correlation between two values of h is defined as

$$\text{cor}(h_i, h_j) = \exp\left\{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\rho}\right\},$$

which allows us to see that values of h near each other will have a higher correlation and thus similar values. This is also why the resulting function is smooth.

Connection to mixed models

It is useful to make connections between this definition of kernel machine regression and mixed models. @liu_semiparametric_2007 demonstrated this by representing $h(\mathbf{x})$ as following a Gaussian process probability distribution,

$$h(\mathbf{x}) \sim \mathcal{GP}(\mathbf{0}, \tau k(\mathbf{x}, \mathbf{x}')),$$

with covariance function k , where \mathbf{x} is a vector of the predictor values, and \mathbf{x}' contains the predictor values of another subject. A Gaussian process is a collection of random variables, of which any finite number follow a multivariate normal distribution [Schulz_2018]. Here, we assume that the expected value of the h function with input \mathbf{x} is $\mathbf{0}$. We use k for the covariance function, which represents the dependence between the function values with inputs \mathbf{x} and \mathbf{x}' : $k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(h(\mathbf{x}) - \mathbf{0})(h(\mathbf{x}') - \mathbf{0})]$.

Now, we can represent h as a collection of variables from a Gaussian process. h follows a multivariate normal distribution,

$$h(\mathbf{x}) \sim N(\mathbf{0}, \tau \mathbf{K}),$$

where $h(\mathbf{x}) = [h(\mathbf{x}_1), h(\mathbf{x}_2), \dots, h(\mathbf{x}_n)]^\top$ and \mathbf{K} is the kernel matrix. The kernel matrix is an $n \times n$ matrix with (i, j) th element $k(\mathbf{x}_i, \mathbf{x}_j)$. Now, returning back to the regression view, we can think of each Y_i as following the distribution,

$$Y_i \stackrel{\text{ind}}{\sim} N(h(\mathbf{x}_i), \sigma^2) \text{ for } i = 1, \dots, n,$$

where σ^2 comes from the variance of the residuals. Here, h can be interpreted as a random effect.

Toy example

In the following section, we illustrate kernel machine regression with a toy example.

Consider the following case where we want to model the relationship between a single predictor and a response variable. Suppose the true relationship between x and Y is defined $Y = e^{\frac{x}{10}} + 2 \sin(\frac{x}{2})$. We simulate 51 equally spaced observations of x from 0 to 25, with error $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, 0.25)$.

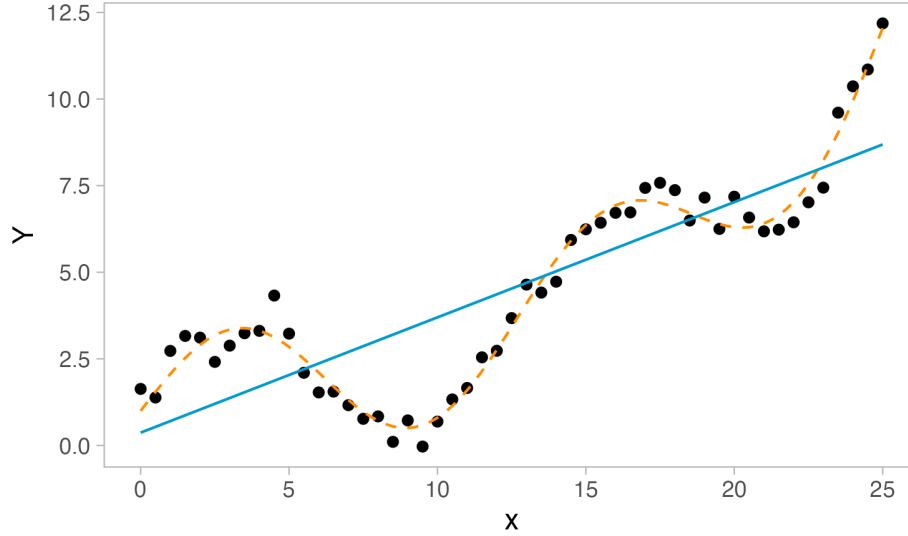


Figure 1: Non-linear data with a true distribution (orange) and a fitted linear regression (blue).

Figure @ref(fig:toy1) illustrates the shape of our simulated non-linear data and the fit proposed by a simple linear regression. We can observe that the linear regression fails to capture the true non-linear relationship. In this case, this would lead to an underestimation of the true association between x and Y . Now, we will try to capture this relationship using kernel machine regression.

To visualize how kernel machine regression works as a moving weighted average, we can consider a query point of 12.5.

Figure @ref(fig:toy2) identifies the query point and assigns corresponding weights to the neighboring points based on a normal distribution, which shares the same density as the Gaussian kernel. In this case, we will specify $\rho = 2$, which is synonymous with assigning the weights using a normal distribution with $\sigma^2 = 1$. We can see how an appropriate estimate for $h(12.5)$ can be obtained by taking a weighted average of the Y 's, with those observations nearby weighted the most heavily.

Now, we fit a kernel machine regression on this data with $\rho = 2$ using the `stats` package in R [R_core_team_r_2013].

We can see in Figure @ref(fig:toy3) that kernel machine regression captures the complex non-linear relationship between Y and x and closely follows the true distribution. We do note, though, that the estimation is less precise at the tails, where there is less information provided by local observations. We can also use this example to consider the effect of various values of ρ on the smoothness of the h function.

Figure @ref(fig:toyrho) demonstrates the effect of relatively smaller and larger values of ρ on h . Decreasing the value of ρ allows kernel machine regression to overfit to the noise in the data by relaxing the dependence of neighboring values of h to each other. On the other hand, increasing the value of ρ enforces more dependence in h and as such results in an underfit estimation. Hence, the choice of ρ has a strong effect on the performance of kernel machine regression.

Variable selection

Now that we have defined kernel machine regression, we can extend it to the Bayesian paradigm. @bobb_bayesian_2015 showed that the Bayesian approach can outperform frequentist kernel machine regression because simultaneous variable selection and estimation can better capture the exposure-response relationship. In this section, we discuss the two methods for Bayesian variable selection in BKMR: hierarchical variable selection and component-wise variable selection [bobb_bayesian_2015].

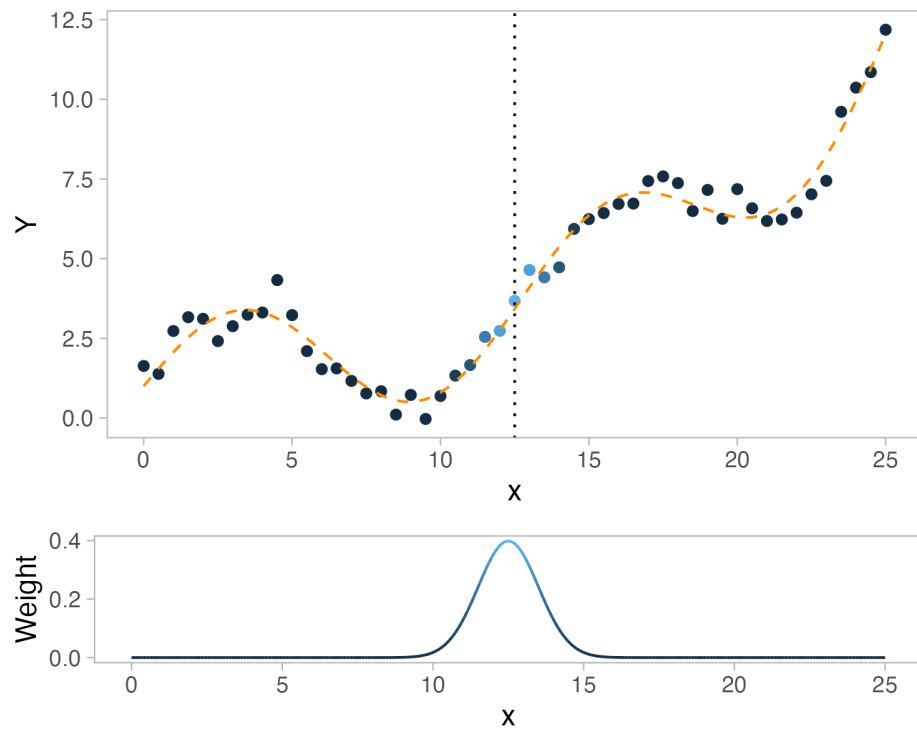


Figure 2: A query point of 12.5 and the weights of neighboring observations based on a Gaussian kernel

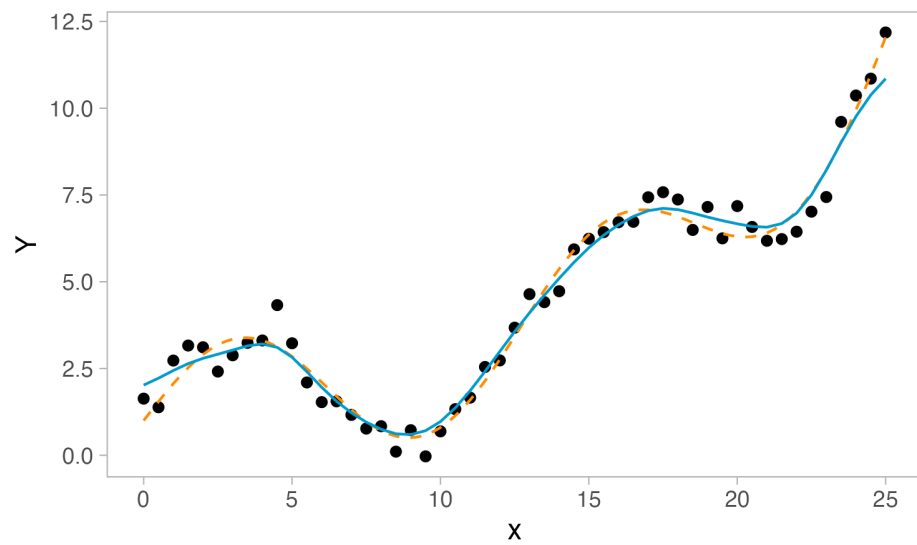


Figure 3: Fitted kernel machine regression with $\rho = 2$.

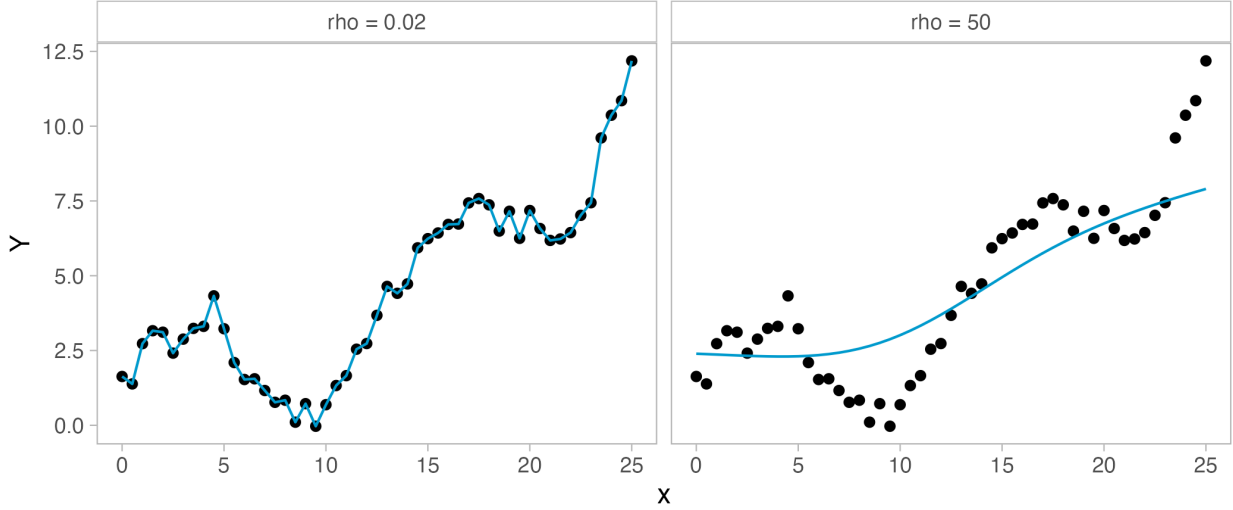


Figure 4: Fitted kernel machine regression with $\rho = 0.02$ and $\rho = 50$.

Component-wise selection follows the same framework as variable selection in a typical Bayesian multiple regression except, instead of augmenting each predictor, we augment the kernel function as

$$k(\mathbf{x}, \mathbf{x}' | \mathbf{r}) = \exp \left\{ - \sum_{m=1}^M r_m (x_m - x'_m)^2 \right\},$$

where $\mathbf{r} = [r_1, \dots, r_M]^\top$. We define $r_m = 1/\rho_m$, the inverse of the tuning parameter ρ_m for each \mathbf{x}_m , as we want r_m to be able to take on the value of 0 when a variable is removed during the selection process. We now define the kernel matrix $\mathbf{K}_{\mathbf{X}, \mathbf{r}}$ as the $n \times n$ matrix with (i, j) th element $k(\mathbf{x}_i, \mathbf{x}_j | \mathbf{r})$. To allow r_m to equal 0 with non-zero probability, we first define an indicator variable determining whether or not a predictor is included in the variable, which is distributed as

$$\delta_m \sim \text{Bernoulli}(\pi),$$

where π is the prior probability of inclusion. Now, we can assume a “slab-and-spike” prior on r_m , distributed as

$$r_m | \delta_m \sim \delta_m f(r_m) + (1 - \delta_m) P_0,$$

where $f(\cdot)$ is some pdf with support \mathbb{R}^+ , and P_0 denotes the density with point mass at 0.

While this process of component-wise variable selection works well in a typical multiple regression setting, it can lead to unreliable estimates in situations where the predictors are highly correlated with each other, which is common in exposure mixture studies. In this case, the correlated components contribute similar information to the model, and component-wise variable selection is not able to distinguish which predictor is important. BKMR deals with this problem by introducing hierarchical variable selection.

Hierarchical variable selection involves partitioning the predictors $\mathbf{x}_1, \dots, \mathbf{x}_M$ into groups \mathcal{S}_g with $g = 1, \dots, G$. These groups should be selected by the user based on a prior knowledge, with the aim of keeping within-group correlation high and between-group correlation low. The indicators from $r_m | \delta_m$ are now distributed as

$$\begin{aligned}\delta_{\mathcal{S}_g}|\omega_g &\sim \text{Multinomial}(\omega_g, \boldsymbol{\pi}_{\mathcal{S}_g}), g = 1, \dots, G, \\ \omega_g &\sim \text{Bernoulli}(\boldsymbol{\pi}),\end{aligned}$$

where $\delta_{\mathcal{S}_g} = \{\delta_m | \mathbf{x}_m \in \mathcal{S}_g\}$ and $\boldsymbol{\pi}_{\mathcal{S}_g}$ are vectors of indicator variables and prior probabilities, respectively, of a predictor \mathbf{x}_m in group \mathcal{S}_g entering the model. By this approach, at most one predictor in each group is allowed to enter the model.

While hierarchical variable selection resolves the issue of multicollinearity, it requires specifying subgroups of predictors a priori and assumes that one predictor in each group can capture the information of the rest. Hence, care should be taken to justify the partitioning of predictors when taking this approach.

Note also that the posterior means of δ_m generated from these variable selection procedures represent the posterior probability of inclusion of \mathbf{x}_m . We can interpret these posterior “inclusion probabilities” as measures of the relative importance of each predictor. These measures can be used to understand the contribution of each exposure or covariate to the health outcome of interest in the model.

Prior specification

In this section, we specify the default prior distributions and parameters used by the BKMR algorithm [bobb_bayesian_2015].

BKMR, by default, assumes $\rho_m = 1/r_m \sim \text{Unif}(a_r, b_r)$, a flat prior between a_r and b_r for which the default values are 0 and 100, respectively [bobb_introduction_2017]. This defines the prior probability of ρ as equally distributed across any value from 0 to 100. This inverse of this prior corresponds to the slab component of the “slab-and-spike” prior, where $r_m|\delta_m \sim \delta_m \text{Unif}^{-1}(a_r, b_r) + (1 - \delta_m)P_0$. As a flat prior, this distribution should be chosen when we have no prior knowledge about the smoothness of the exposure-response function, with hyperparameters a_r and b_r selected to represent the range of values we expect ρ to potentially span.

We have seen that the smoothness of a kernel machine regression responds strongly to different values of $r_m = 1/\rho$, and, accordingly, the model fit of BKMR is sensitive to their prior distribution. In general, the posterior inclusion probabilities generated from the variable selection procedure are particularly sensitive to this prior, though their relative importance tends to remain stable [bobb_statistical_2018]. As such, the BKMR algorithm also offers the options to define uniform and gamma priors for the $r_m = 1/\rho$.

Moreover, BKMR assumes that the prior probability of including a predictor (δ_m) or group of predictors (ω_g) in the model is distributed $\pi \sim \text{Beta}(a_\pi, b_\pi)$. The default hyperparameters are $a_\pi = b_\pi = 1$, which represent a flat, uninformative prior between 0 and 1. When the hierarchical selection approach is applied, equal values for $\boldsymbol{\pi}_{\mathcal{S}_g}$, the probabilities of inclusion for each component in group \mathcal{S}_g , are assumed.

Finally, BKMR assumes that the inverse of the variance of the residuals is distributed $\sigma^{-2} \sim \text{Gamma}(a_\sigma, b_\sigma)$, with default values of $a_\sigma = b_\sigma = 0.001$, and that the vertical scale of h is parameterized by $\lambda \equiv \tau\sigma^{-2} \sim \text{Gamma}(a_\lambda, b_\lambda)$, with default values of a_λ, b_λ such that the mean and variance of λ are both equal to 10.

The MCMC algorithm

Briefly, we discuss the algorithm used in the BKMR package [bobb_bayesian_2015; bobb_statistical_2018], with commentary on its implications for the model fitting process.

BKMR uses a Markov chain Monte Carlo (MCMC) algorithm with a mix of Gibbs and Metropolis-Hastings samplers to estimate the posterior distributions. In particular, a Gibbs step is used to update the distribution of σ^2 while a Metropolis-Hastings step is used to update the distribution of λ . For component-wise and hierarchical variable selection, $(\mathbf{r}, \boldsymbol{\delta}, \boldsymbol{\omega})$ are sampled jointly using a Metropolis-Hastings sampling scheme.

While each distribution generated by the Gibbs step is always accepted, the distributions for λ and r_m generated by the Metropolis-Hastings steps are accepted based on an acceptance rate [wagaman_probability_2021]. The standard deviation of the proposal distribution controls the acceptance rate and as such acts as a tuning parameter [bobb_example_2017]. In general, increasing the standard deviation leads to lower acceptance rates. Acceptance rates that are too low lead to slower convergence, but rates that are too high can cause convergence to a non-optimal distribution.

A major computational limitation of BKMR is that at each iteration of the MCMC algorithm, the $n \times n$ augmented kernel matrix $\mathbf{K}_{\mathbf{Z},\mathbf{r}}$ must be inverted multiple times. BKMR can employ a Gaussian predictive process which involves specifying a set of l points, or “knots,” that are a subset of the predictor space. The vector of predictors can be approximated by projection on this lower dimensional space, which allows the algorithm to perform inversions on an $l \times l$ matrix. A general suggestion is to use this approach to speed up the algorithm when n is large and to specify $l \approx n/10$ [bellavia_statistical_2021].

Bayesian semiparametric regression

differences w/ bkmr

- makes distributional assumptions about the dataset
- kernel regression is computationally intensive w/ large datasets but can handle many predictors
- bsr highly dependent on the choice of function

Defining notation for spline regression:

- \mathbf{x}

BSR specific notation:

- \mathbf{x}

Spline regression

We begin by introducing spline regression, with attention to its specific implementation in BSR. Spline regression is a semiparametric regression technique that can be used to capture non-linear effects. In this introduction, we follow the presentation of spline regression provided by antonelli_estimating_2020, with additional details and explanation from hastie_elements_2009.

Spline regression defines the regression relationship as

$$Y_i = f(\mathbf{x}_i) + \epsilon_i,$$

where f is defined by a set of basis functions and ϵ_i is a random variable from $\epsilon \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$. BSR uses natural spline bases. In order to understand how these are constructed, we start with a broad definition of basis expansions, before exploring linear, cubic, and then natural spline bases.

We determine the basis expansion by considering a piece-wise function of X_i with some order P and some set of K knots defining $K + 1$ disjoint intervals. BSR places knots at uniformly sized quantiles within the boundaries of X_i . The most commonly used orders are $P = 1, 2,$ and 4 , the constant, linear, and cubic splines, respectively. To begin, let us consider a continuous piece-wise linear spline basis (i.e., $P = 2$) of a one-dimensional X with two interior knots. In this case, we use the following four basis functions:

$$b_1(X) = 1, \quad b_2(X) = X, \quad b_3(X) = (X - \xi_1)_+, \quad b_4(X) = (X - \xi_2)_+,$$

where ξ_1 and ξ_2 are the two interior knots, and t_+ denotes the positive part. These bases are used to construct the regression model $f(X) = \sum_{j=1}^4 \beta_j b_j(X)$, which requires estimating $K+P = 4$ parameters. We can check the continuity restrictions at the knots by seeing that $f(\xi_1^-) = \beta_1 + \xi_1 \beta_2$ and $f(\xi_1^+) = \beta_1 + \xi_1 \beta_2 + (\xi_1 - \xi_1) \beta_3$ are equal.

Now, in the case of exposure mixtures, we want smoother functions that can capture the non-linear relationship between the response and the predictors. We can achieve this by increasing the order to $P = 4$ and using a cubic spline, with continuous first and second derivatives at the knots. The cubic spline is the lowest-order spline for which knot-discontinuity cannot be detected by the human eye. For example, for one X with two interior knots, we use the following six basis functions:

$$\begin{aligned} b_1(X) &= 1, & b_2(X) &= X, & b_3(X) &= X^2, \\ b_4(X) &= X^3, & b_5(X) &= (X - \xi_1)_+^3, & b_6(X) &= (X - \xi_2)_+^3. \end{aligned}$$

Now, the regression model is defined as $f(X) = \sum_{j=1}^6 \beta_j b_j(X)$ and requires estimating $K+P = 6$ parameters. It can be shown that $f'(\xi_i^-) = f'(\xi_i^+)$ and $f''(\xi_i^-) = f''(\xi_i^+)$.

However, the behavior of polynomials near the boundaries of X , where there is less information, can be erratic. Natural cubic splines, also referred to as just natural splines, address this by imposing an additional restriction of linearity at the boundaries of X . Paradoxically, this also leads to a simpler model with four fewer parameters to estimate. A general definition of the K basis functions for a natural spline with interior knots $\xi_j, j = 1, \dots, K$, is given by:

$$\begin{aligned} N_1(X) &= 1, & N_2(X) &= X, & N_{k+2}(X) &= d_k(X) - d_{K-1}(X), \\ d_k(X) &= \frac{(X - \xi_k)_+^3 - (X - \xi_K)_+^3}{\xi_K - \xi_k}. \end{aligned}$$

Here, the regression model is defined as $f(X) = \sum_{j=1}^K \beta_j N_j(X)$, with K parameters. BSR uses natural splines to specify the regression relationship.

Toy example

In the following section, we illustrate spline regression using the same toy example as presented in Section [bkmr toy example]. See Section [bkmr toy example] and Figure @ref(fig:toy1) for details on the parameters used to generate simulated data.

As in Section @ref(bkmr_toy), we consider a case where we want to model the relationship between a single predictor and a response variable, where the true relationship between x and Y is defined $Y = e^{\frac{x}{10}} + 2 \sin(\frac{x}{2})$. We fit a series of linear, cubic, and then natural spline regressions to illustrate the general framework of a spline regression.

Figure @ref(fig:toy4) illustrates the fit proposed by a linear spline regression with order $P = 2$. We can see that the implementation of knots allows for even a linear fit to capture more of the nuances in this nonlinear relationship, as compared to a standard linear regression.

However, this linear spline regression is still unable to fully estimate the nonlinearity in our example. Increasing the order to $P = 4$ with a cubic spline regression offers additional flexibility. Figure toy5 illustrates the fit proposed by this model. Here, we can see the benefits of using a cubic polynomial relationship in a nonlinear setting: the estimated relationship is continuous at the knots (specifically, the second derivatives are 0), and the nonlinear relationship has been flexibly captured.

Our final modification involves imposing linearity constraints

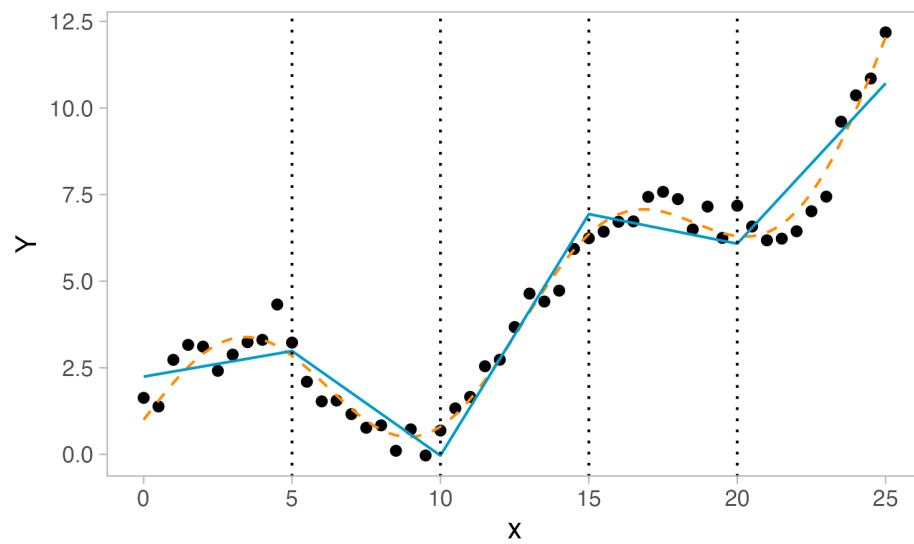


Figure 5: Linear spline regression.

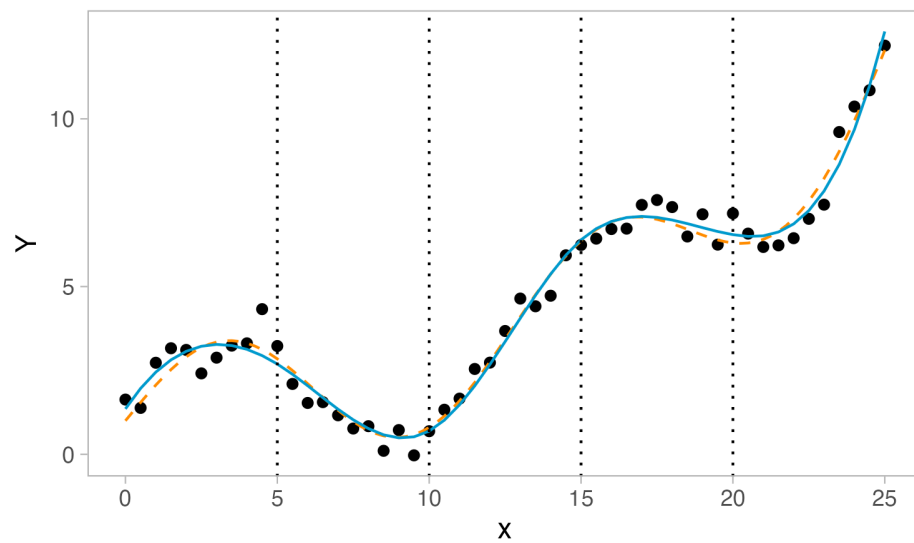


Figure 6: Cubic spline regression.

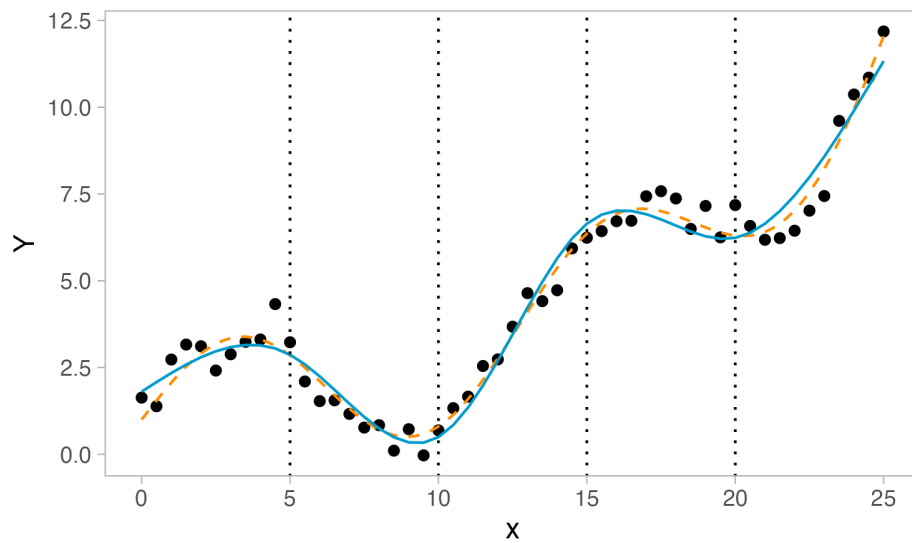


Figure 7: Natural spline regression.

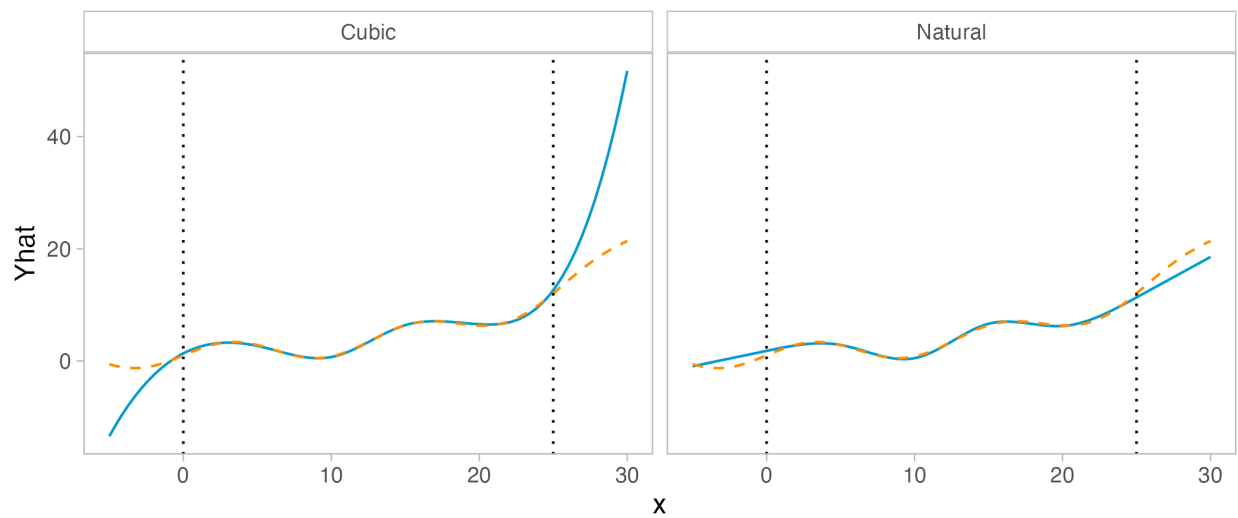


Figure 8: Natural and cubic spline regression extrapolated outside the bounds of x .

Model formulation in BSR

Now that we have defined natural splines, we introduce BSR, following the presentation in @antonelli_estimating_2020. We demonstrate the construction of f in BSR by first assuming a two-dimensional case with predictors X_1 and X_2 . We define

$$\begin{aligned} f(X_1) &= \tilde{X}_1 \beta_1, & f(X_2) &= \tilde{X}_2 \beta_2, \\ f(X_1, X_2) &= \tilde{X}_1 \beta_1 + \tilde{X}_2 \beta_2 + \tilde{X}_{12} \beta_{12}, \end{aligned}$$

where $\tilde{X}_m = [b_{m1}(X_m), \dots, b_{md}(X_m)]$ represents a d -dimensional basis function expansion for $m = 1, 2$, and $\tilde{X}_{12} = [b_{11}(X_1)b_{21}(X_2), b_{11}(X_1)b_{22}(X_2), \dots, b_{1d}(X_1)b_{2d}(X_2)]$ represents a d^2 -dimensional basis expansion of the interaction between X_1 and X_2 . Note that we must explicitly model the effect of the interaction term by assuming a multiplicative interaction between the basis functions of the predictors.

Extending to the multi-dimensional setting, BSR assumes the following general model formulation:

$$\begin{aligned} f(\mathbf{x}_i) &= \sum_{h=1}^p f^{(h)}(\mathbf{x}_i), \\ f^{(h)}(\mathbf{x}_i) &= \sum_{j_1=1}^M \tilde{x}_{ij_1} \beta_{j_1}^{(h)} + \sum_{j_1=2}^M \sum_{j_2 < j_1} \tilde{x}_{ij_1 j_2} \beta_{j_1 j_2}^{(h)} + \dots, \end{aligned}$$

where $f^{(h)}(\mathbf{x}_i)$ includes a summation of all p -way interactions. The inclusion of all p -way interactions makes the model far too overparameterized. Moreover, $f(\mathbf{x}_i)$ is a sum of k different functions $f^{(h)}(\mathbf{x}_i)$ where a value for k is selected in order to capture all exposure effects in the model. Each of the k functions has the same functional form, and so the regression coefficients for a function $f^{(h)}(\mathbf{x}_i)$ are only identifiable up to a constant — this means that there are multiple sets of coefficients that could be estimated from the same data.

Sparsity inducing priors

In order to handle the overparameterization and non-identifiability of the model, BSR implements multivariate sparsity inducing priors. In this section, we follow the presentation provided in @antonelli_estimating_2020.

First, we define indicators $\zeta = \{\zeta_{jh}\}$ representing whether the j th exposure is included in the h th function:

$$\begin{aligned} P(\zeta_{jh}) &= \tau_h^{\zeta_{jh}} (1 - \tau_h)^{1 - \zeta_{jh}} I(A_h \not\subset A_{h'} \forall h' \neq h \text{ or } A_h = \{\}), \\ &\text{where } A_h = \{j : \zeta_{jh} = 1\}. \end{aligned}$$

Here, the indicators follow a Bernoulli distribution with prior probability of inclusion τ_h . The posterior means of ζ can be interpreted as measures of relative variable importance. We include an indicator function $I()$ that represents whether the function h contains a unique set of predictors. This indicator ensures that no function contains predictors that are a subset of those in another function, in which case this function would be redundant and thus removed from the model entirely.

Now, we assume a multivariate slab-and-spike prior on the regression coefficients:

$$P(\beta_X^{(h)}|\zeta) = \left(1 - \prod_{j \in S} \zeta_{jh}\right) P_0 + \left(\prod_{j \in S} \zeta_{jh}\right) \psi_1(\beta_S^{(h)}),$$

where S is some subset of $1, 2, \dots, p$.

Here, P_0 denotes the density with point mass at $\mathbf{0}$, and $\psi_1()$ is a multivariate normal distribution with mean $\mathbf{0}$ and covariance Σ_β , a diagonal matrix with $\sigma^2 \sigma_\beta^2$ on the diagonals.

Prior specification

- hyper prior vs. empirical bayes strategy for sigma beta
- lower bound on slab variance
- properties of priors for higher-order interactions

The MCMC algorithm

- include brief detail on how MCMC algorithm deals with lower-order interactions that are a subset of higher-order interactions

Detecting interactions

Discuss options for detecting interactions. Briefly describe MLR .

In Section @ref(motivation), we highlighted the challenges of analytically testing for the presence of interactions in exposure mixture studies. These challenges motivated a theoretical exploration of BKMR and BSR in Sections @ref(bkmr) and @ref(bsr). Now, we discuss and compare the options that BKMR and BSR provide for inference on the presence of interactions. We also include discussion on theoretical explanations for why we might prefer one over the other.

BKMR

Since the flexible h function in kernel machine regression allows us to forgo any assumptions about the nature of the relationship between the health outcome and predictors, BKMR can potentially capture complex interactions between predictors. The challenge with using BKMR to do this, however, is that there is no formal framework for conducting inference on the presence of interactions.

And, if we specify hierarchical variable selection, then

How to detect interactions:

- can visualize relationship at various quantiles of other predictor
- can conduct inference on inter-quantile difference — has yet to be a formal assessment of how this approach performs in various settings

BSR

BSR can conduct inference