

P-Spline Vector Generalized Additive Models and Its Applications

Chanatda Somchit, Chris Wild, and Thomas Yee

Department of Statistics
University of Auckland

1 December 2015

- 1 Introduction to VGLMs and VGAMs
- 2 P-spline VGAMs
- 3 Simulation study
- 4 Examples
- 5 References

VGLMs and VGAMs are the extension of class of **GLMs and GAMs** to include a class of multivariate regression models [Yee and Wild, 1996].

- **VGLMs** model each parameter as a linear combination of the covariates,

$$\eta_j(\mathbf{x}) = \beta_j^T \mathbf{x} = \sum_{k=1}^p \beta_{(j)k} x_k, \quad j = 1, \dots, M,$$

- **VGAMs** extend VGLMs to

$$\eta_j(\mathbf{x}) = \sum_{k=1}^p f_{(j)k}(x_k), \quad j = 1, \dots, M.$$

The current class of VGLMs/VGAMs is very large and includes many statistical distributions and models:

- univariate and multivariate distributions,
- categorical data analysis,
- quantile and expectile regression,
- time series, survival analysis,
- extreme value analysis, nonlinear regression,
- reduced-rank regression, ordination, etc.

The underlying algorithm of VGAMs is

- Modified vector backfitting using vector splines.
- But...it is not easy to integrate the automatic numerical procedure used to determine the shape of non-linear terms from the data.

We aim to...

- integrate an automatic procedure for estimating the smoothing parameters to the VGAM framework.

To achieve this.....

- the ideas of GAMs based on penalized regression splines proposed by Marx and Eilers [1998] and Wood [2006] are generalized to the VGAM class.

Therefore, ...

- we develop VGAMs based on penalized regression splines with P-spline smoothers, which we call 'P-spline VGAMs'.

As a result, ...

- P-spline VGAMs can be transformed into the VGLM framework,
- and maximized by penalized iteratively re-weighted least squares (P-IRLS).
- The computational procedure for the automatic and stable multiple smoothing parameter selection can be implemented.
- The issue of determining the shape of the smooth terms can be resolved.

- The underlying ideas of P-splines are
 - to use B-splines as basis functions, and a large number of equally-spaced knots are used,
 - but to prevent the problem of overfitting, a discrete approximate wiggleness penalty is applied to the model fitting objective,

$$\sum_{s=1}^S \left(\Delta^{[d]} a_s \right)^2 = \mathbf{a}^T \mathbf{P}_{[d]} \mathbf{a}, \quad (1)$$

where $\mathbf{P}_{[d]} = \mathbf{D}_{[d]}^T \mathbf{D}_{[d]}$, $\mathbf{D}_{[d]}$ is the matrix consisting of d^{th} order difference of the coefficients, and \mathbf{a} is a parameter vector.

- Here is a linear combination of B-spline basis functions:

$$f(x_i) = \sum_{s=1}^S a_s B_s(x_i). \quad (2)$$

The basic structure of P-spline VGAMs is

$$g_j(\theta_j) = \eta_j(\mathbf{x}) = f_{(j)1}(x_1) + \cdots + f_{(j)p}(x_p), \quad j = 1, \dots, M, \quad (3)$$

- where $f_{(j)k}$ are represented using B-splines, and are centered for identifiability.
- If there is an intercept, $x_1 = 1$,
- All P-spline VGAM smooth components are estimated simultaneously.

Allow the P-spline VGAM approach to be more general for use in many more situations

In practice, we may wish to constrain the effects of a single covariate

- to be the same for different η_j , or
- to have no effect for others.
- For example

$$\begin{aligned}\eta_1 &= \beta_{(1)1} + f_{(1)2}(x_2) + f_{(1)3}(x_3) \\ \eta_2 &= \beta_{(2)1} + f_{(1)2}(x_2).\end{aligned}\tag{4}$$

- **Yee and Wild [1996]** introduced ‘**constraint matrices**’ applied directly to the linear/additive predictors to control how the covariates act – “**constraints on the functions**”.

These constraints

- are very useful for categorical models, such as, a bivariate odds-ratio model, and the proportional odds model,
- lead the VGLM/VGAM approach to be more general for use in most situations.

We will generalize these ideas to the P-spline VGAM framework.

In general for P-spline VGAMs, we represent the models as

$$\begin{aligned}\eta &= \mathbf{f}_1(x_1) + \cdots + \mathbf{f}_p(x_p) \\ &= \mathbf{H}_1 \mathbf{f}_1^*(x_1) + \cdots + \mathbf{H}_p \mathbf{f}_p^*(x_p),\end{aligned}\tag{5}$$

where

- $\mathbf{H}_1, \dots, \mathbf{H}_p$ are known full column-rank ‘constraint matrices’,
- $\mathbf{f}_k^* = (f_{(1)k}(x_k), \dots, f_{(R_k)k}(x_k))^T$ is a vector consisting of a possibly reduced set of smooth functions.
- Each smooth term is centered for identifiability.
- No constraints at all, $\mathbf{H}_k = \mathbf{I}_M$.

The smooth predictor vector $\boldsymbol{\eta}$ can be now written as

$$\boldsymbol{\eta}(\mathbf{x}_i) = \sum_{k=1}^p \mathbf{H}_k \mathbf{X}_{ik}^* \boldsymbol{\beta}_k^*, \quad (6)$$

where

- $\mathbf{X}_{ik}^* = \mathbf{x}_{ik}^{*T} \otimes \mathbf{I}_{R_k}$, $\mathbf{x}_{ik}^* = (B_{k:1}(x_{ik}), \dots, B_{k:S_k}(x_{ik}))^T$ is a vector of B-splines generated at the values of x_k and i th observation,
- $\boldsymbol{\beta}_k^* = \text{vec} \begin{pmatrix} \boldsymbol{\beta}_{(1)k}^{*T} \\ \vdots \\ \boldsymbol{\beta}_{(R_k)k}^{*T} \end{pmatrix}$, where $\boldsymbol{\beta}_{(r_k)k}^* = (a_{(r_k)k:1}, \dots, a_{(r_k)k:S_k})^T$, for $r_k = 1, \dots, R_k$.

Setting up P-spline VGAMs as penalized VGLMs (3)

- Let $\beta^* = (\beta_1^{*T}, \dots, \beta_p^{*T})^T$ be a vector containing all of the possibly reduced sets of B-spline coefficients in the models,
- $\mathbf{X}_{\text{vam}} = \left((\mathbf{X}_{\text{am}} \tilde{\mathbf{e}}_1) \otimes \mathbf{H}_1 \mid (\mathbf{X}_{\text{am}} \tilde{\mathbf{e}}_2) \otimes \mathbf{H}_2 \mid \dots \mid (\mathbf{X}_{\text{am}} \tilde{\mathbf{e}}_p) \otimes \mathbf{H}_p \right)$,
where
 - \mathbf{X}_{vam} is the model matrix for P-spline VGAMs,
 - \mathbf{X}_{am} is the 'additive model' model matrix for one η_j , and
 - $\tilde{\mathbf{e}}_k = \begin{pmatrix} \mathbf{0} \\ \mathbf{I}_{(s_k \times s_k)} \\ \mathbf{0} \end{pmatrix}$.

- Then,

$$\boldsymbol{\eta} = \mathbf{X}_{\text{vam}} \boldsymbol{\beta}^*. \quad (7)$$

- Equation (7) is just the form of VGLMs, therefore

$$\ell(\boldsymbol{\beta}^*) = \sum_{i=1}^n w_i \ell\{\eta_1(\mathbf{x}_i), \dots, \eta_M(\mathbf{x}_i)\}. \quad (8)$$

Now, it is a good chance to control the model's smoothness by adding **a wiggleness penalty** to the log-likelihood objective of $\ell(\beta^*)$.

- The penalty term for P-spline VGAMs is given by

$$\begin{aligned}
 J(\lambda) &= \sum_{k=1}^p \sum_{j=1}^{R_k} \lambda_{(j)k} \beta_{(j)k}^T \mathbf{D}_{[d]k}^T \mathbf{D}_{[d]k} \beta_{(j)k} \\
 &= \sum_{k=1}^p \beta_k^{*T} \left\{ \left(\mathbf{D}_{[d]k}^T \mathbf{D}_{[d]k} \right) \otimes \text{diag} \left(\lambda_{(1)k}, \dots, \lambda_{(R_k)k} \right) \right\} \beta_k^* \\
 &= \sum_{k=1}^p \beta_k^{*T} \mathbf{P}_{\lambda_k}^* \beta_k^*,
 \end{aligned}$$

where $\mathbf{P}_{\lambda_k}^* = \left(\mathbf{D}_{[d]k}^T \mathbf{D}_{[d]k} \right) \otimes \text{diag} \left(\lambda_{(1)k}, \dots, \lambda_{(R_k)k} \right)$.

- Then, the quadratic penalty on the parameter vector β^* for P-spline VGAMs is given by

$$J(\lambda) = \beta^{*T} \mathbf{P}_\lambda^* \beta^*,$$

- where $\mathbf{P}_\lambda^* = \text{blockdiag}(\mathbf{P}_{\lambda_1}^*, \dots, \mathbf{P}_{\lambda_p}^*)$.
- The penalized log-likelihood for P-spline VGAMs is then

$$\ell^*(\beta^*) = \ell(\beta^*) - \frac{1}{2} \beta^{*T} \mathbf{P}_\lambda^* \beta^*. \quad (9)$$

- **Newton-Raphson algorithm** is applied for maximizing the log-likelihood (9),

$$\beta^{*(t+1)} = \beta^{*(t)} + \mathcal{I}(\beta^{*(t)})^{-1} U(\beta^{*(t)}). \quad (10)$$

- **Maximizing** $\ell^*(\beta^*)$ using (10) leads to an iterative solution for penalized iteratively reweighted least squares (P-IRLS) of

$$\beta^{*(t+1)} = \left(\mathbf{X}_{\text{vam}}^T \mathbf{W}^{(t)} \mathbf{X}_{\text{vam}} + \mathbf{P}_\lambda^* \right)^{-1} \left(\mathbf{X}_{\text{vam}}^T \mathbf{W}^{(t)} \mathbf{z}^{(t)} \right). \quad (11)$$

- Given values for λ , $\hat{\beta}^{*(t+1)}$ is the solution to

$$\min_{\beta^*} \quad \left(\mathbf{z} - \mathbf{X}_{\text{vam}} \beta^* \right)^T \mathbf{W} \left(\mathbf{z} - \mathbf{X}_{\text{vam}} \beta^* \right) + \beta^{*T} \mathbf{P}_\lambda^* \beta^*, \quad (12)$$

where

- $\mathbf{W} = \text{blockdiag}(\mathbf{W}_1, \dots, \mathbf{W}_n)$,
 - Fisher scoring: $(\mathbf{W}_i)_{jk} = -w_i E \left(\frac{\partial^2 \ell_i}{\partial \eta_j \partial \eta_k} \right)$,
- $\mathbf{z} = (\mathbf{z}_1^T, \dots, \mathbf{z}_n^T)^T$ and $\mathbf{u} = (\mathbf{u}_1^T, \dots, \mathbf{u}_n^T)^T$, where
 $(\mathbf{u}_i)_j = w_i \frac{\partial \ell_i}{\partial \eta_j}$, and $\mathbf{z}_i = \boldsymbol{\eta}_i + (\mathbf{W}_i)^{-1} \mathbf{u}_i$.

- **Data augmentation** is applied to the adjusted dependent vector, regressors, and weights:

$$\mathbf{z}' = \begin{pmatrix} \mathbf{z} \\ \vartheta \end{pmatrix}, \quad \mathbf{X}'_{\text{vam}} = \begin{pmatrix} \mathbf{X}_{\text{vam}} \\ \tilde{\mathbf{P}}_{\lambda}^* \end{pmatrix}, \quad \mathbf{W}' = \text{blockdiag}(\mathbf{W}, \mathbf{I}_{\vartheta}),$$

where $\mathbf{P}_{\lambda}^* = \tilde{\mathbf{P}}_{\lambda}^{T*} \tilde{\mathbf{P}}_{\lambda}^*$, $\vartheta = \sum_{k=1}^p (S_k - d) \cdot M$.

- Therefore, (12) can be replaced by the equivalent:

$$\min_{\beta^*} \left(\mathbf{z}' - \mathbf{X}'_{\text{vam}} \beta^* \right)^T \mathbf{W}' \left(\mathbf{z}' - \mathbf{X}'_{\text{vam}} \beta^* \right). \quad (13)$$

- We convert the GLS system of equations to OLS:

$$\mathbf{z}''(t) = \mathbf{X}''_{\text{vam}}(t) \beta^* + \varepsilon''(t). \quad (14)$$

- **In P-IRLS**, we fit OLS model above using the data augmentation of \mathbf{z} , \mathbf{W} , and \mathbf{X}_{vam} until the convergence is achieved.

- Given an estimate for β^* , the multiple smoothing parameter selection for penalized least squares above can be solved by the minimization of the GCV or the 'unbiased risk estimator' (UBRE) w.r.t. the multiple smoothing parameters.
- The UBRE score for the P-spline VGAM approach:

$$\nu_u(\lambda) = \frac{1}{nM} \left\| \sqrt{\mathbf{W}}(z - \mathbf{X}_{\text{VAM}}\beta^*) \right\|^2 - 1 + \frac{2}{nM} \text{tr}(\mathbf{A}_\lambda).$$

- $\mathbf{A}_\lambda = \mathbf{X}_{\text{VAM}}^T \mathbf{W} \mathbf{X}_{\text{VAM}} \left(\mathbf{X}_{\text{VAM}}^T \mathbf{W} \mathbf{X}_{\text{VAM}} + \mathbf{P}_\lambda \right)^{-1}$ is the hat matrix.
- λ enters the UBRE score through \mathbf{A}_λ .
- $\text{tr}(\mathbf{A}_\lambda)$ represents the estimated effective degree of freedom (EDF).

Estimating smoothing parameters (2)

- each working penalized linear model of the P-IRLS iteration, $\nu_u(\lambda)$ is minimized w.r.t. λ (**performance iteration**).
- The two steps:

Step 1 Obtain an estimate of β^* via:

$$\beta^{*(t+1)} = \underset{\beta^*}{\operatorname{argmax}} \ell^*(\beta^*).$$

Step 2 Obtain an estimate of λ via:

$$\lambda^{(t+1)} = \underset{\lambda}{\operatorname{argmin}} \nu_u(\lambda).$$

The two steps of estimations above are iterated until convergence is met.

- Here, we employ **the computational approach developed by Wood [2004]** to minimize the UBRE or GCV scores.

- The simulation study was based on the model:

$$y_{i1}^* = \beta_{(1)1} + \beta_{(1)2} x_{i2} + f_{(1)3}(x_{i3}) + f_{(1)4}(x_{i4}) + \varepsilon_{i1},$$

$$y_{i2}^* = \beta_{(2)1} + \beta_{(2)2} x_{i2} + f_{(2)3}(x_{i3}) + f_{(2)4}(x_{i4}) + \varepsilon_{i2},$$

- The binary responses y_{i1} and y_{i2} are determined according to the rule:

$$\begin{cases} y_{ij} = 1 & \text{if } y_{ij}^* > 0 \\ y_{ij} = 0 & \text{if } y_{ij}^* \leq 0 \end{cases} ; j = 1, 2.$$

- The three test functions, $f_{(1)3}(x_{i3}) = \cos(2\pi x_{i3})$, $f_{(2)3}(x_{i3}) = 2 \sin(\pi x_{i4})$ and $f_{(1)4}(x_{i4}) = f_{(2)4}(x_{i4}) = 0 x_{i4}$.
- Three uniform covariates on $(0, 1)$ were simulated.
- The error terms $(\varepsilon_{i1}, \varepsilon_{i2})$:

$$\begin{pmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \end{pmatrix} \stackrel{iid}{\sim} N \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \right).$$

Simulation Study (2): a semiparametric bivariate probit model

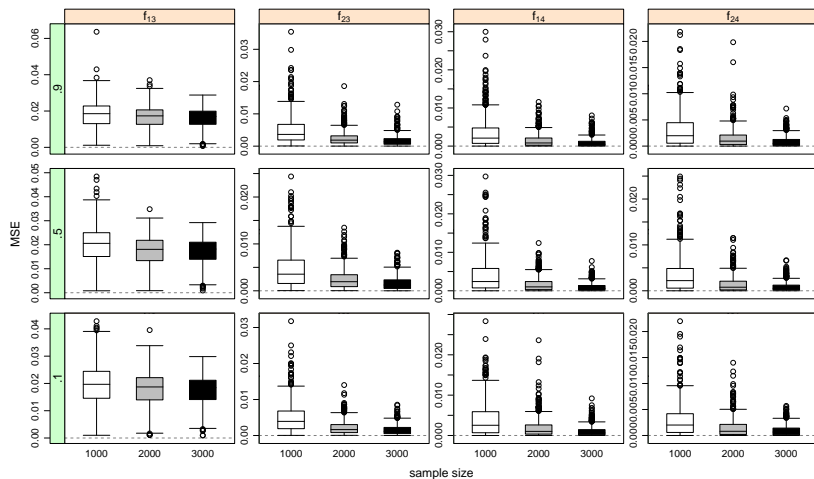


Figure: Boxplots of the MSE of $\hat{f}_{(1)3}$, $\hat{f}_{(2)3}$, $\hat{f}_{(1)4}$, and $\hat{f}_{(2)4}$, when employing the P-spline VGAM approach. The numbers .1, .5, .9 in the y-axis captions denote the three different correlations, ρ .

Simulation Study (3): a semiparametric bivariate probit model

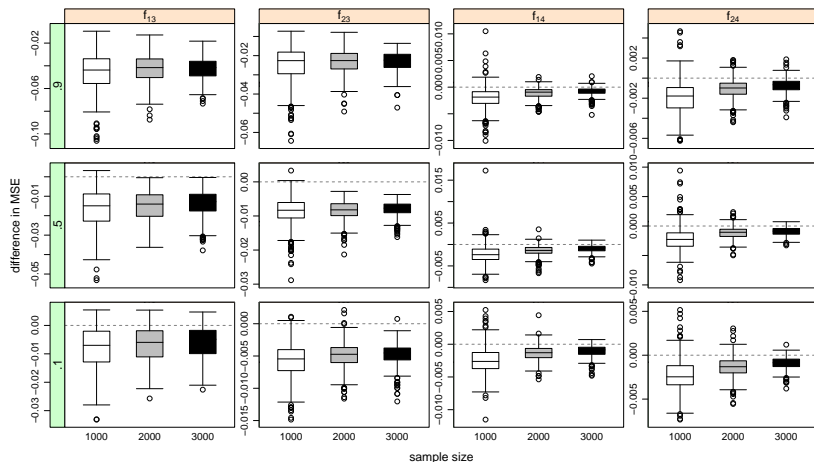


Figure: Boxplots of the difference in MSE between the P-spline VGAM approach and the VGAM approach of $\hat{f}_{(1)3}$, $\hat{f}_{(2)3}$, $\hat{f}_{(1)4}$, and $\hat{f}_{(2)4}$. Negative values indicate that the new approach performs better than the alternative.

Examples: Mackerel egg survey (1)

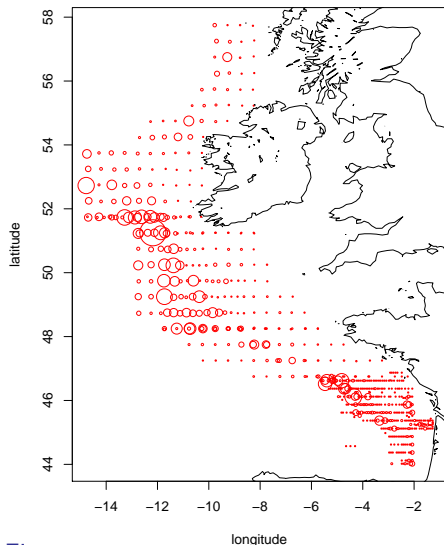


Figure: Observed mackerel eggs densities per square metre of sea surface as assessed by net samples.

- The data from 1992 mackerel egg survey
- The response of interest is the egg counts (`egg.count`).
- The covariates of interest are:
 - latitude (`lat`),
 - longitude (`lon`)
 - water temperature at a depth of 20m (`temp.20`),
 - the sea bed depth at the sampling location (`b.dept`),
 - distance from 200m sea bed contour (`c.dist`).

***Borchers et al. [1997] and Wood [2006]

Examples: Mackerel egg survey (2)

An additive Poisson model is fitted to the count response (egg.count), with a mean given by

$$E[\text{egg.count}_i] = g_i \times [\text{net area}]_i,$$

where g_i is the density of eggs, per square metre of sea surface, at i^{th} sampling location.

$$\log(E[\text{egg.count}_i]) = f_i + \log([\text{net area}]_i),$$

where $f_i = \log(g_i)$.

```
fit.ps1 <- psvgam(egg.count ~ ps(lat, 5) + ps(lon, 5) +  
                  ps(temp.20m, 5) + ps(b.depth, 5) +  
                  ps(c.dist, 5) + offset(log.net.area),  
                  data = mack1, family = poissonff)
```

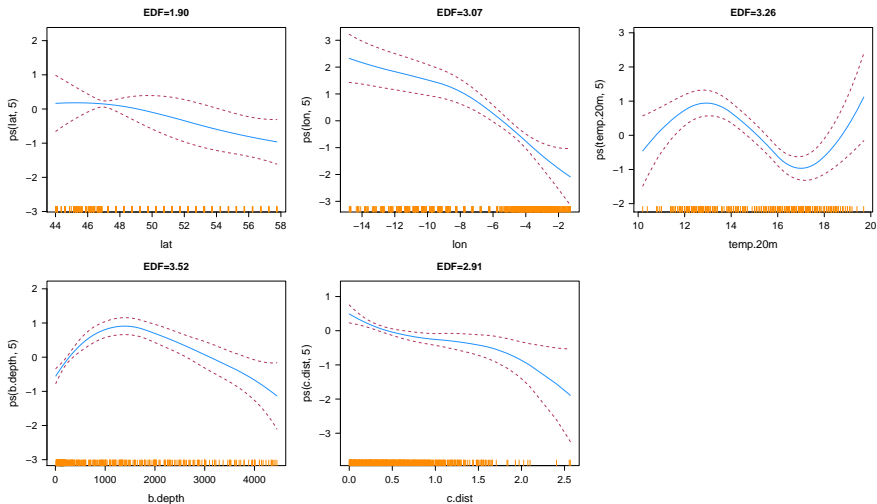
- The dispersion parameter estimate indicates the existence of overdispersion in this data set.
- Trying the negative binomial distribution:

```
fit.ps2 <- psvgam(egg.count ~ ps(lat, 5) + ps(lon, 5) +  
                 ps(temp.20m, 5) + ps(b.depth, 5) +  
                 ps(c.dist, 5) + offset(log.net.area),  
                 data = mack1, family = negbinomial)
```

- The negative binomial model is more appropriate than the Poisson model.

Examples: Mackerel egg survey (4)

Estimated smooth terms for the mackerel model `fit.ps2`.



Examples: Mackerel egg survey (5)

Model Predictions

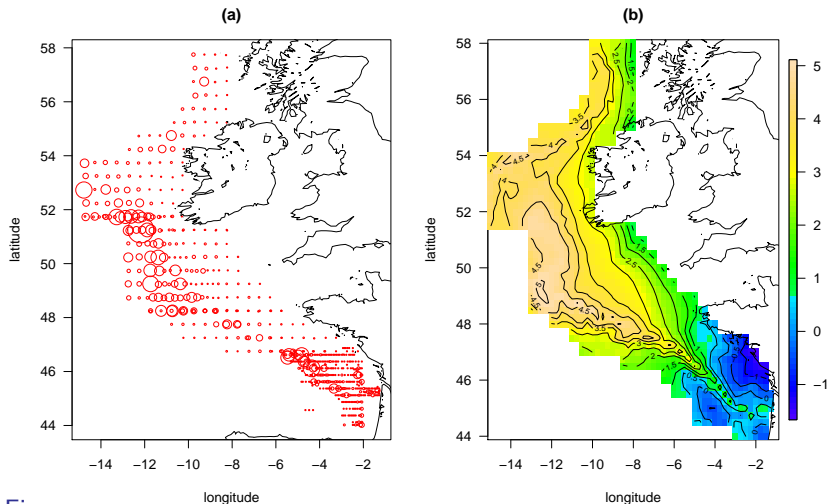


Figure: (a) Observed mackerel eggs densities per square metre of sea surface as assessed by net samples. (b) Predicted log densities of mackerel eggs over the survey area, according to the model fit.ps1.

Let's investigate how the probability of having a household cat and a household dog is related to people's ages.

- We fit a **nonparametric bivariate logistic model** to the example of cat and dog pet ownership presented by Yee [2015].
- For homogeneity, we restrict the analysis to a subset of 2569 European women and remove any missing values.

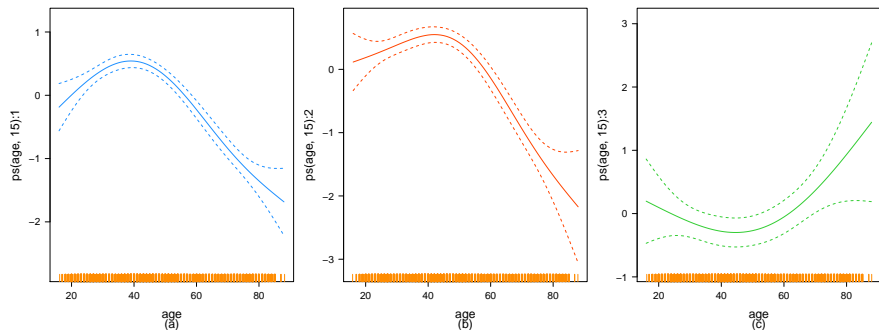
$$\begin{aligned}\eta_1 &= \text{logit } P(Y_1 = 1|x_2) = \beta_{(1)1} + f_{(1)2}(x_2), \\ \eta_2 &= \text{logit } P(Y_2 = 1|x_2) = \beta_{(2)1} + f_{(2)2}(x_2), \\ \eta_3 &= \log \psi = \beta_{(3)1} + f_{(3)2}(x_2),\end{aligned}\tag{15}$$

where ψ is the odds ratio.

```
fitps.cd1 <- psvgam(cbind(cat, dog) ~ ps(age, 15),  
                   binom2.or(zero = NULL),  
                   data = women.eth0.catdog)
```

Examples: Cat and dog data (2)

Plots of fitted component functions, according to the model `fitps.cd1`.



Examples: Cat and dog data (3)

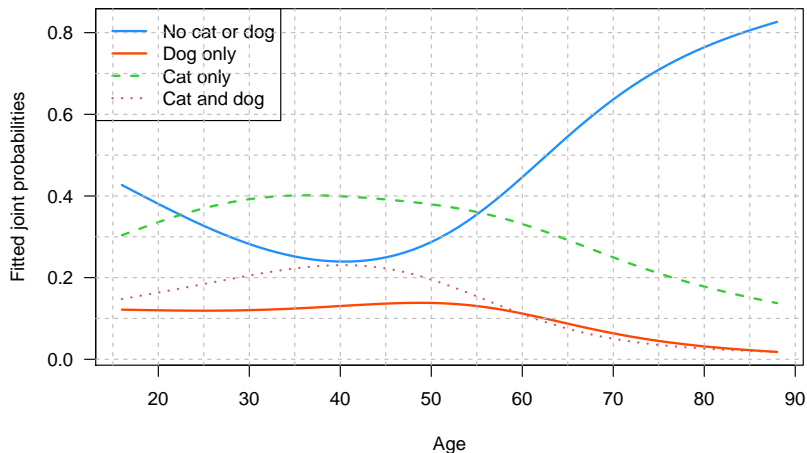


Figure: Estimated probabilities for all four combinations (both cats and dogs, cats only, dogs only, and no cats or dogs) of a subset of European women using the P-spline VGAM approach.

- D. Borchers, S. Buckland, I. Priede, and S. Ahmadi. Improving the precision of the daily egg production method using generalized additive models. *Canadian Journal of Fisheries and Aquatic Sciences*, 54(12):2727–2742, 1997.
- B. D. Marx and P. H. Eilers. Direct generalized additive modeling with penalized likelihood. *Computational Statistics & Data Analysis*, 28(2):193–209, 1998.
- S. N. Wood. Stable and efficient multiple smoothing parameter estimation for generalized additive models. *Journal of the American Statistical Association*, 99:673–686, 2004.
- S. N. Wood. *Generalized additive models: An introduction with R*. Chapman & Hall/CRC, Boca Raton, FL, USA, 2006.
- T. W. Yee. *Vector Generalized Linear and Additive Models: With an Implementation in R*. Springer, New York, NY, USA, 2015. To appear.
- T. W. Yee and C. J. Wild. Vector generalized additive models. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58:481–493, 1996.