



# ***Additive Smooth Models: basic theory***

Simon N. Wood

University of Bath

# Additive smooth models

- ⑥ Consider a univariate response  $y$  and corresponding predictors  $x_1, x_2, x_3 \dots$  (possibly vectors).
- ⑥ An additive smooth model has a structure like

$$y_i = \mathbf{X}_i \boldsymbol{\beta} + f_1(x_{1i}) + f_2(x_{2i}) + \dots + \epsilon_i$$

— the  $f_j$  are **smooth functions** and  $\mathbf{X}_i \boldsymbol{\beta}$  is the linear predictor for any *strictly parametric model components*.

- ⑥ The  $\epsilon_i$  are 0 mean r.v.s, with variance  $\sigma^2$ . Normality is assumed if CIs or  $H_0$  tests are required.
- ⑥ The additive smooth structure offers a nice balance of flexibility and structure.

# ASM representation

- How can the  $f_j$  be estimated? Use a basis expansion for each smooth term. The model becomes

$$y_i = \mathbf{X}_i \boldsymbol{\beta} + \sum_{k=1}^{K_1} \beta_{1k}^* b_{1k}(x_{1i}) + \sum_{k=1}^{K_2} \beta_{2k}^* b_{2k}(x_{2i}) + \cdots + \epsilon_i$$

- Ignorance of parametric form of  $f_j$  OK, if  $K_j$  large enough and *basis functions*,  $b_{jk}(x)$ , chosen carefully.
- Only **unknowns** are now parameters. Absorbing  $b_{jk}(x_{ji})$  into  $\mathbf{X}_i$  and  $\beta_{jk}^*$  into  $\boldsymbol{\beta}$ , the model becomes

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (\text{a linear model!})$$

# Not over-fitting

- ⑥ To avoid model mis-specification bias we made  $K_j$ s fairly large (not too large or computational efficiency is lost).  $\Rightarrow$  There is a danger of overfitting.
- ⑥ To avoid overfitting penalize lack of smoothness. e.g. estimate  $\beta$  by minimizing not just RSS, but

$$\mathcal{S}(\beta) = \text{RSS} + \sum_j \lambda_j \times [\text{wiggleness of } f_j]$$

- ⑥ Smoothing parameters  $\lambda_j$ , assumed known for moment, while e.g.  $[\text{wiggleness of } f_j] = \int f_j''(x)^2 dx$ .

# Model estimation

- Given basis functions, most wiggleness measures (including  $\int f_j''(x)^2 dx$ ) can be written as  $\beta^T \mathbf{S}_j \beta$  where  $\mathbf{S}_j$  contains known coefficients.
- $f_j$  only identifiable to within an additive constant  $\Rightarrow$  impose sum-to-zero constraints, and absorb into basis.
- Fitting objective becomes, minimize

$$\mathcal{S}(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2 + \sum \lambda_j \beta^T \mathbf{S}_j \beta \text{ w.r.t. } \beta.$$

- Resulting estimate is  $\hat{\beta} = \mathbf{B}\mathbf{y}$  where  $\mathbf{B} = (\mathbf{X}^T \mathbf{X} + \sum_j \mathbf{S}_j)^{-1} \mathbf{X}^T$  (not for computational use!)

## Inference based on $\hat{\beta}$

- ⑥ Covariance matrix of  $\hat{\beta}$  easily derived from that of  $\mathbf{y}$

$$\mathbf{V}_{\hat{\beta}} = \mathbf{B}\mathbf{B}^T \sigma^2.$$

- ⑥ If we also assume normality (i.e.  $\mathbf{y} \sim N(\mathbf{0}, \mathbf{I}\sigma^2)$ ) then

$$\hat{\beta} \sim N(\mathbb{E}(\hat{\beta}), \mathbf{B}\mathbf{B}^T \sigma^2) \quad (1)$$

- ⑥ But  $\mathbb{E}(\hat{\beta}) \neq \beta$  unless  $\beta = 0$  (strictly *in null space of penalty*)... so (1) sometimes useful for hypothesis testing, but probably not otherwise.

# Bayesian Inference

- ⑥ We penalize wiggleness because we think that a smooth truth is 'more probable' than a wiggly one.
- ⑥ So let 'model wiggleness' have a negative exponential prior distribution (Wahba, '83; Silverman '85),

$$\pi(\boldsymbol{\beta}) \propto e^{-\frac{1}{2} \sum \lambda_j \boldsymbol{\beta}^\top \mathbf{S}_j \boldsymbol{\beta}}$$

- ⑥ The model says  $\mathbf{y}|\boldsymbol{\beta} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{I}\sigma^2)$ , so Bayes rule

$$\Rightarrow \boldsymbol{\beta}|\mathbf{y} \sim N(\hat{\boldsymbol{\beta}}, (\mathbf{X}^\top \mathbf{X} + \sum_j \mathbf{S}_j)^{-1} \sigma^2)$$

- ⑥ Credible intervals for any function of  $\boldsymbol{\beta}$  follow.

## ***Inference based on fit***

- ⑥ There's always an unpenalized model with fitted value behaviour that is 'close' to a penalized model with given smoothing parameters.
- ⑥ Hence, conditional on smoothing parameters, the distributions of F-ratio statistics for comparison of penalized models should be approximately the distribution of the equivalent un-penalized analogues.
- ⑥ So F-ratio testing can be used for model comparison, provided we can work out the approximate 'equivalent degrees of freedom' of the penalized models.



# Degrees of freedom

- ⑥ The point of penalizing the fit is to reduce the model's freedom to vary. How many degrees of freedom does the penalized fit have?
- ⑥ Without penalization,  $\tilde{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ .
- ⑥ With penalization,  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X} + \sum_j \lambda_j \mathbf{S}_j)^{-1} \mathbf{X}^\top \mathbf{y}$ .
- ⑥ So  $\hat{\beta} = \mathbf{F} \tilde{\beta}$  where  $\mathbf{F} = (\mathbf{X}^\top \mathbf{X} + \sum_j \lambda_j \mathbf{S}_j)^{-1} \mathbf{X}^\top \mathbf{X}$ .
- ⑥  $\mathbf{F}_{ii} = \partial \hat{\beta}_i / \partial \tilde{\beta}_i$  is a measure of the DoF of the  $i^{\text{th}}$  parameter, while  $\text{tr}(\mathbf{F}) \equiv \text{the whole model DoF}$ .
- ⑥  $\hat{\sigma}^2 = \|\mathbf{y} - \mathbf{X} \hat{\beta}\|^2 / (n - \text{tr}(\mathbf{F}))$ .

# Smoothness selection

- ⑥ So far everything is conditional on  $\lambda$ : how should  $\lambda$  values be chosen?
- ⑥ There are two main choices:
  1. Use prediction error criteria, such as cross validation, GCV, AIC, BIC etc.
  2. Use the Bayes smoothing model to decompose each smooth into fixed effect and random effect components: the model can then be estimated by maximum likelihood or REML, with  $\lambda$  treated as variance parameters.

`mgcv::gam` does 1, while `mgcv::gamm` does 2.  
Concentrate on 1, here.

# Ordinary cross validation

- ⑥ Leave one out cross validation seeks to minimize estimated mean square prediction error to select  $\lambda$ .

$$\text{minimize } V_o(\boldsymbol{\lambda}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\mu}_i^{[-i]})^2 \text{ w.r.t. } \boldsymbol{\lambda}$$

where  $\hat{\mu}_i^{[-i]}$  is predicted  $y_i$  from fit to all data except  $y_i$ .

- ⑥ If  $\mathbf{A} = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \sum_j \lambda_j \mathbf{S}_j)^{-1} \mathbf{X}^\top$ , so  $\hat{\boldsymbol{\mu}} = \mathbf{A}\mathbf{y}$ , then

$$V_g(\boldsymbol{\lambda}) = \frac{1}{n} \sum_{i=1}^n \frac{(y_i - \hat{\mu}_i)^2}{(1 - A_{ii})^2}$$

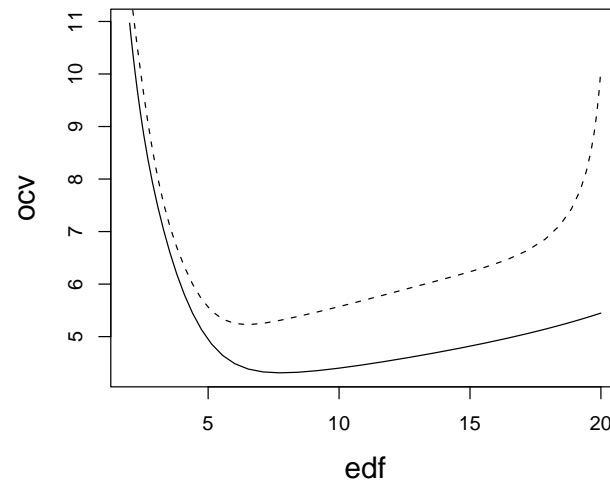
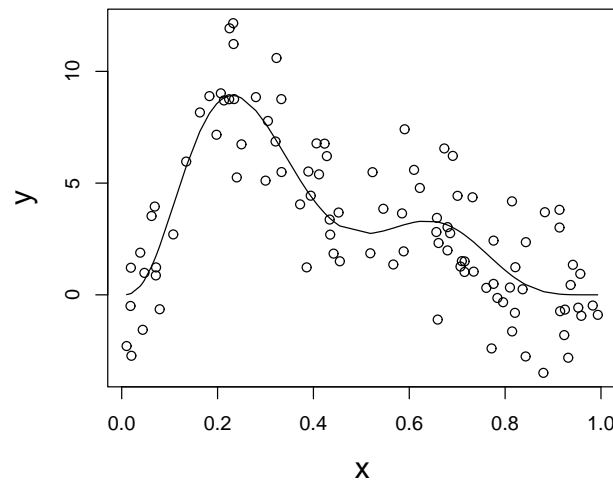
# OCV not invariant

- For any  $\perp$  matrix  $\mathbf{Q}$ , fitting objective

$$\mathcal{S}_Q(\boldsymbol{\beta}) = \|\mathbf{Q}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\|^2 + \sum \lambda_j \boldsymbol{\beta}^\top \mathbf{S}_j \boldsymbol{\beta}$$

yields identical inferences about  $\boldsymbol{\beta}$  as  $\mathcal{S}(\boldsymbol{\beta})$ .

- But  $\mathcal{S}_Q$  and  $\mathcal{S}$  yield different OCV functions!



# Generalized cross validation

- ⑥ If we choose  $Q$  so that each rotated observation counts equally in the OCV function, then we get GCV (e.g. Craven and Wahba, '79).
- ⑥ The resulting GCV function,

$$V_g(\boldsymbol{\lambda}) = \frac{n \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{[n - \text{tr}(\mathbf{A})]^2}$$

is invariant. Note that  $\text{tr}(\mathbf{A}) = \text{tr}(\mathbf{F})$ .

- ⑥ Main problem solved by `mgcv::gam` is to minimize functions like  $V_g$  efficiently w.r.t.  $\boldsymbol{\lambda}$ .

## Other $\lambda$ selection criteria

- ⑥ Can increase 'cost' per DoF in GCV function

$$V_g(\boldsymbol{\lambda}) = \frac{n \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{[n - \gamma \text{tr}(\mathbf{F})]^2} \quad \gamma \geq 1$$

- ⑥ If  $\sigma^2$  is known, can estimate  $\mathbb{E}(\|\boldsymbol{\mu} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2)/\sigma^2$  by

$$V_u(\boldsymbol{\lambda}) = \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2/\sigma^2 + 2\text{tr}(\mathbf{F})$$

which is basically Mallows' statistic, Craven and Wahba's UBRE or AIC. The penalty per DoF can be increased to give, e.g. BIC.

# Generalized Additive Models

- ⑥ The model is generalized to

$$g\{\mathbb{E}(y_i)\} = \mathbf{X}_i\boldsymbol{\beta} + f_1(x_{1i}) + f_2(x_{2i}) + \dots$$

- ⑥  $g$  is a known 'link function'.
- ⑥  $y_i$  are independent and either
  1. assume  $y_i$  follow an exponential family distribution (use maximum likelihood), or
  2. assume only  $\text{var}(y_i) = \phi V(\mathbb{E}(y_i))$ , where  $V$  is a known function (use quasi-likelihood).

# GAM representation

- As for ASM, basis expansion allows model to be written as a GLM  $g\{\mathbb{E}(y_i)\} = \mathbf{X}_i\boldsymbol{\beta}$ .
- Fit is measured by log-likelihood ( $l(\boldsymbol{\beta})$  or equivalently  $l(\boldsymbol{\mu})$ ), or more conveniently *Deviance*

$$D(\boldsymbol{\mu}) = 2\{l(\mathbf{y}) - l(\boldsymbol{\mu})\}$$

which takes on the role of the RSS.

- Again penalization avoids overfitting.  $\hat{\boldsymbol{\beta}}$

$$\text{minimizes } D(\boldsymbol{\beta}) + \sum \lambda_j \boldsymbol{\beta}^\top \mathbf{S}_j \boldsymbol{\beta} \text{ w.r.t. } \boldsymbol{\beta}$$



# GAM estimation

- ⑥ Penalized MLE is performed by Penalized IRLS.

- ⑥ Let  $\hat{\beta}$  and  $\hat{\mu}$  be current best estimates. Define

$$z_i = g'(\hat{\mu}_i)(y_i - \hat{\mu}_i) + \mathbf{X}_i \hat{\beta} \quad \text{and} \quad W_{ii} = \{V(\hat{\mu}_i)g'(\hat{\mu}_i)^2\}^{-1}$$

- ⑥ The  $\beta$  minimizing

$$\|\sqrt{\mathbf{W}}(\mathbf{z} - \mathbf{X}\beta)\|^2 + \sum \lambda_j \beta^T \mathbf{S}_j \beta$$

is (almost always) an improved  $\hat{\beta}$  estimate.

- ⑥ Iteratively repeating the steps finds the penalized likelihood maximizing estimates.

# GAM Inference

- ⑥ The ASM results have large sample equivalents for GAMs.
- ⑥ If  $\mathbf{B} = (\mathbf{X}^T \mathbf{W} \mathbf{X} + \sum \lambda_j \mathbf{S}_j)^{-1} \mathbf{X}^T \mathbf{W}$  then  $\hat{\boldsymbol{\beta}} \sim N(\mathbb{E}(\hat{\boldsymbol{\beta}}), \mathbf{B} \mathbf{B}^T \phi)$  [frequentist].
- ⑥  $\boldsymbol{\beta} \sim N(\hat{\boldsymbol{\beta}}, (\mathbf{X}^T \mathbf{W} \mathbf{X} + \sum \lambda_j \mathbf{S}_j)^{-1} \phi)$  [Bayesian].
- ⑥ Approximate GLRT and F-ratio tests are possible, as if using a GLM, provided the effective degrees of freedom ( $\text{tr}(\mathbf{F})$ ) is used as the model DoF.

# Smoothness selection

- ⑥ There are two approaches.
- ⑥ **Performance iteration** (Gu '92). Select  $\lambda$  by GCV/UBRE applied to the working linear model at each P-IRLS step.
  - + Fast. - Does not always converge.
- ⑥ **Outer iteration**. Use Deviance,  $D(\hat{\beta})$ , in place of the RSS,  $\|y - X\hat{\beta}\|^2$ , in the GCV or UBRE score definitions, and optimize directly w.r.t.  $\lambda$ .
  - + Good convergence. - Slower.

*Note:* UBRE becomes generalized AIC here.

# Prediction

- ⑥ Because model has a parametric representation as a GLM, prediction from a fitted model is just like prediction from a GLM.
- ⑥ Given new covariate values a 'prediction matrix',  $\mathbf{X}^p$ , is created, exactly as if producing a new model matrix, except that any covariate dependent details of basis function form depend on the *original* fit covariates.
- ⑥ The predicted values are then  $\mu_i^p = \mathbf{X}_i^p \hat{\beta}$ .
- ⑥ By simulating from the posterior of  $\beta$ , samples from the posterior of any function of  $\beta$  can be obtained!

# Open problems

- ⑥ All the distributional results are conditional on the smoothing parameters.
- ⑥ This reduces the reliability of p-values, and causes problems with the component wise performance of credible intervals/confidence intervals.
- ⑥ Work on smoothing parameter unconditional confidence intervals is promising but incomplete (see Wood 2006 *GAMs: An Intro with R*).

## *Other approaches*

There are several R packages providing GAMs or similar models. See the package docs for full references.

- ⑥ `gam` provides Trevor Hasties original backfit GAMs.
- ⑥ `gss` provides full smoothing spline based models (including computationally efficient versions).
- ⑥ `assist` is an alternative for full spline smoothing.
- ⑥ `gamlss` moves beyond smooth models of the mean of a response to model variance skew and kurtosis as smooth functions of covariates.
- ⑥ `vgam` looks at multivariate extensions and much more.