

GAMs: Model Selection

David L Miller, Eric Pedersen, and Gavin L Simpson
August 6th, 2016

Overview

- Model selection
- Shrinkage smooths
- Shrinkage via double penalty (`select = TRUE`)
- Confidence intervals for smooths
- p values
- `anova()`
- AIC

Model selection

Model selection

Model (or variable) selection — and important area of theoretical and applied interest

- In statistics we aim for a balance between *fit* and *parsimony*
- In applied research we seek the set of covariates with strongest effects on y

We seek a subset of covariates that improves *interpretability* and *prediction accuracy*

Shrinkage & additional penalties

Shrinkage & additional penalties

Smoothing parameter estimation allows selection of a wide range of potentially complex functions for smooths...

But, cannot remove a term entirely from the model because the penalties used act only on the *range space* of a spline basis. The *null space* of the basis is unpenalised.

- **Null space** — the basis functions that are smooth (constant, linear)
- **Range space** — the basis functions that are wiggly

Shrinkage & additional penalties

mgcv has two ways to penalize the null space, i.e. to do selection

- *double penalty approach* via `select = TRUE`
- *shrinkage approach* via special bases for thin plate and cubic splines

Other shrinkage/selection approaches are available

Double-penalty shrinkage

\mathbf{S}_j is the smoothing penalty matrix & can be decomposed as

$$\mathbf{S}_j = \mathbf{U}_j \mathbf{\Lambda}_j \mathbf{U}_j^T$$

where \mathbf{U}_j is a matrix of eigenvectors and $\mathbf{\Lambda}_j$ a diagonal matrix of eigenvalues (i.e. this is an eigen decomposition of \mathbf{S}_j).

$\mathbf{\Lambda}_j$ contains some **0**s due to the spline basis null space — no matter how large the penalty λ_j might get no guarantee a smooth term will be suppressed completely.

To solve this we need an extra penalty...

Double-penalty shrinkage

Create a second penalty matrix from \mathbf{U}_j , considering only the matrix of eigenvectors associated with the zero eigenvalues

$$\mathbf{S}_j^* = \mathbf{U}_j^* \mathbf{U}_j^{*T}$$

Now we can fit a GAM with two penalties of the form

$$\lambda_j \boldsymbol{\beta}^T \mathbf{S}_j \boldsymbol{\beta} + \lambda_j^* \boldsymbol{\beta}^T \mathbf{S}_j^* \boldsymbol{\beta}$$

Which implies two sets of penalties need to be estimated.

In practice, add `select = TRUE` to your `gam()` call

Shrinkage

The double penalty approach requires twice as many smoothness parameters to be estimated. An alternative is the shrinkage approach, where \mathbf{S}_j is replaced by

$$\tilde{\mathbf{S}}_j = \mathbf{U}_j \tilde{\mathbf{\Lambda}}_j \mathbf{U}_j^T$$

where $\tilde{\mathbf{\Lambda}}_j$ is as before except the zero eigenvalues are set to some small value ϵ .

This allows the null space terms to be shrunk by the standard smoothing parameters.

Use `s(..., bs = "ts")` or `s(..., bs = "cs")` in **mgcv**

Empirical Bayes...?

\mathbf{S}_j can be viewed as prior precision matrices and λ_j as improper Gaussian priors on the spline coefficients.

The impropriety derives from \mathbf{S}_j not being of full rank (zeroes in Λ_j).

Both the double penalty and shrinkage smooths remove the impropriety from the Gaussian prior

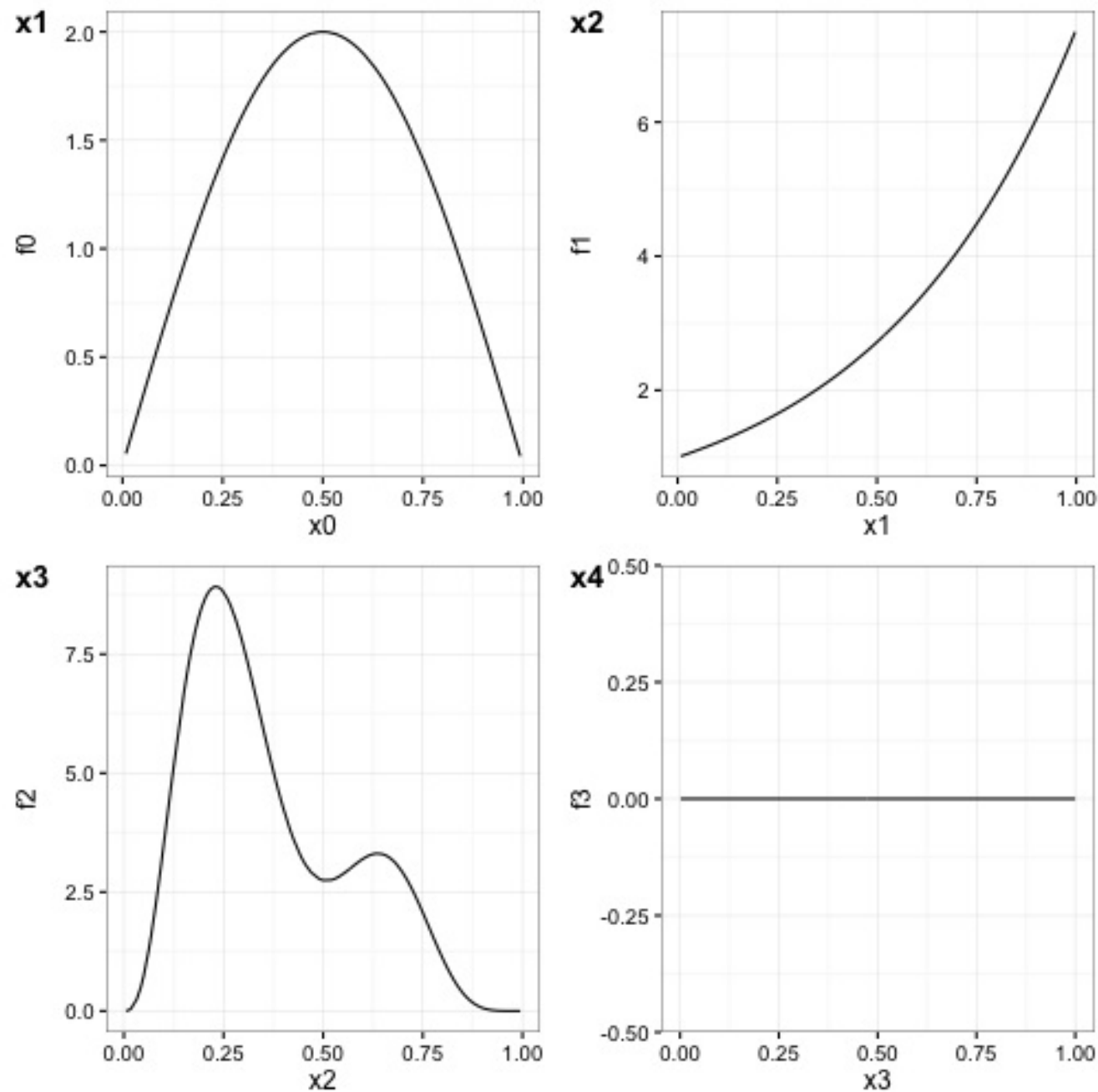
Empirical Bayes...?

- **Double penalty** — makes no assumption as to how much to shrink the null space. This is determined from the data via estimation of λ_j^*
- **Shrinkage smooths** — assumes null space should be shrunk less than the wiggly part

Marra & Wood (2011) show that the double penalty and the shrinkage smooth approaches

- performed significantly better than alternatives in terms of *predictive ability*, and
- performed as well as alternatives in terms of variable selection

Example



- Simulate Poisson counts
- 4 known functions
- 2 spurious covariates

Example

Family: poisson
Link function: log

Formula:
 $y \sim s(x0) + s(x1) + s(x2) + s(x3) + s(x4) + s(x5)$

Parametric coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	1.21758	0.04082	29.83	<2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

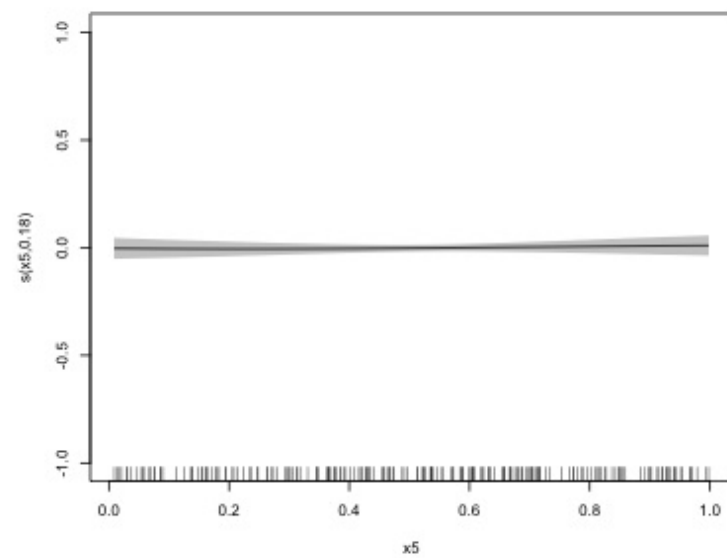
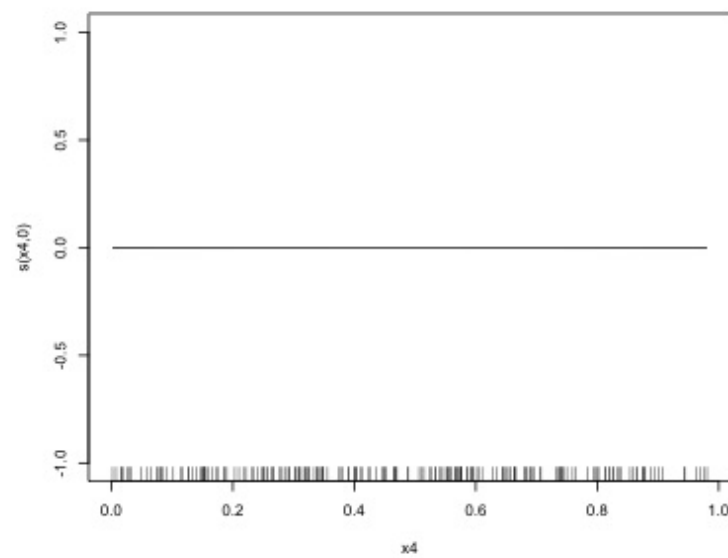
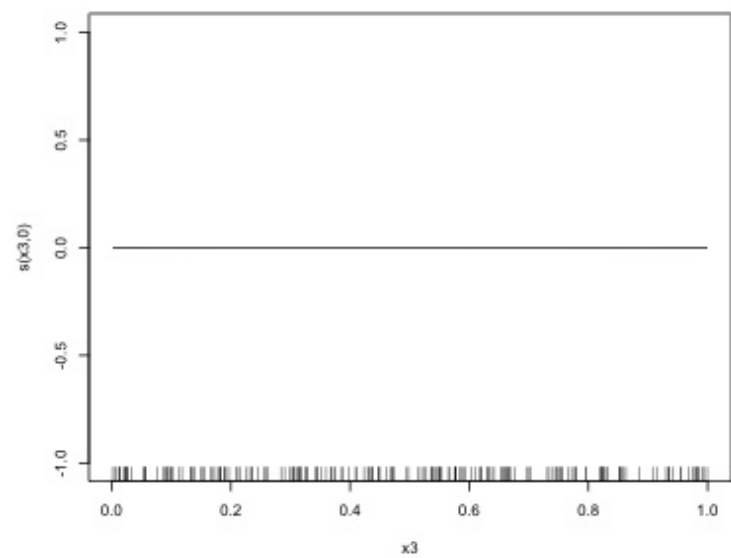
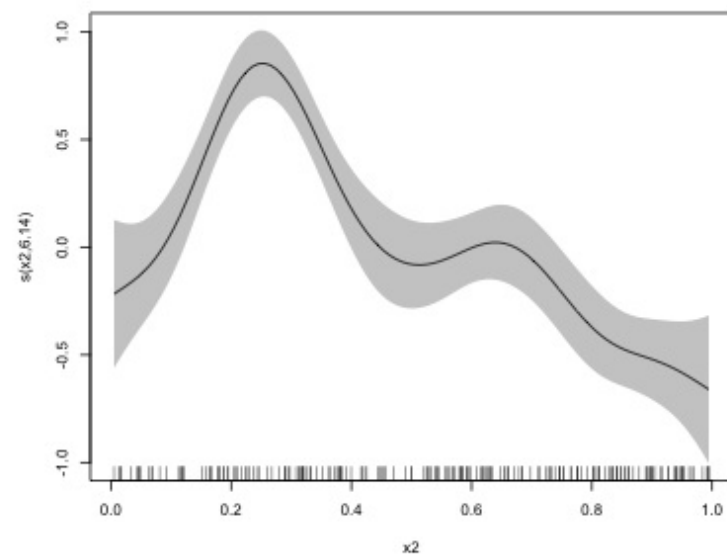
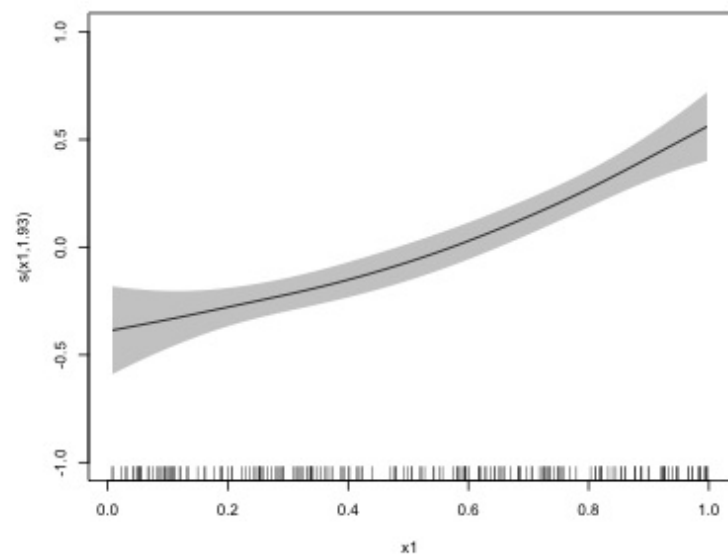
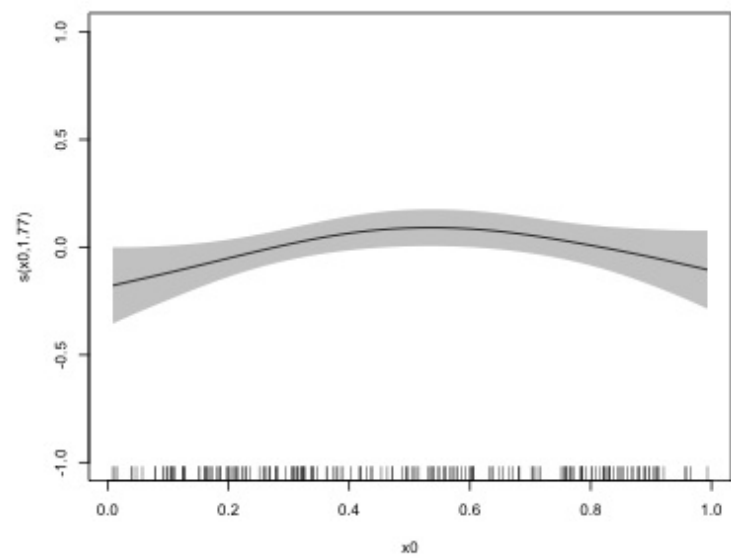
Approximate significance of smooth terms:

	edf	Ref.df	Chi.sq	p-value	
s(x0)	1.7655119	9	5.264	0.0397	*
s(x1)	1.9271039	9	65.356	<2e-16	***
s(x2)	6.1351372	9	156.204	<2e-16	***
s(x3)	0.0002618	9	0.000	0.4088	
s(x4)	0.0002766	9	0.000	1.0000	
s(x5)	0.1757146	9	0.195	0.2963	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

R-sq.(adj) = 0.545 Deviance explained = 51.6%
-REML = 430.78 Scale est. = 1 n = 200

Example



Confidence intervals for smooths

Confidence intervals for smooths

`plot.gam()` produces approximate 95% intervals (at ± 2 SEs)

What do these intervals represent?

Nychka (1988) showed that standard Wahba/Silverman type Bayesian confidence intervals on smooths had good **across-the-function** frequentist coverage properties.

Confidence intervals for smooths

Marra & Wood (2012) extended this theory to the generalised case and explain where the coverage properties failed:

Musn't over-smooth too much, which happens when λ_j are over-estimated

Two situations where this might occur

1. where true effect is almost in the penalty null space,
 $\hat{\lambda}_j \rightarrow \infty$
2. where $\hat{\lambda}_j$ difficult to estimate due to highly correlated covariates
 - if 2 correlated covariates have different amounts of

Don't over-smooth

In summary, we have shown that Bayesian componentwise variable width intervals... for the smooth components of an additive model **should achieve close to nominal *across-the-function coverage probability***, provided only that we do not over-smooth so heavily... Beyond this requirement not to oversmooth too heavily, the results appear to have rather weak dependence on smoothing parameter values, suggesting that the neglect of smoothing parameter variability should not significantly degrade interval performance.

Confidence intervals for smooths

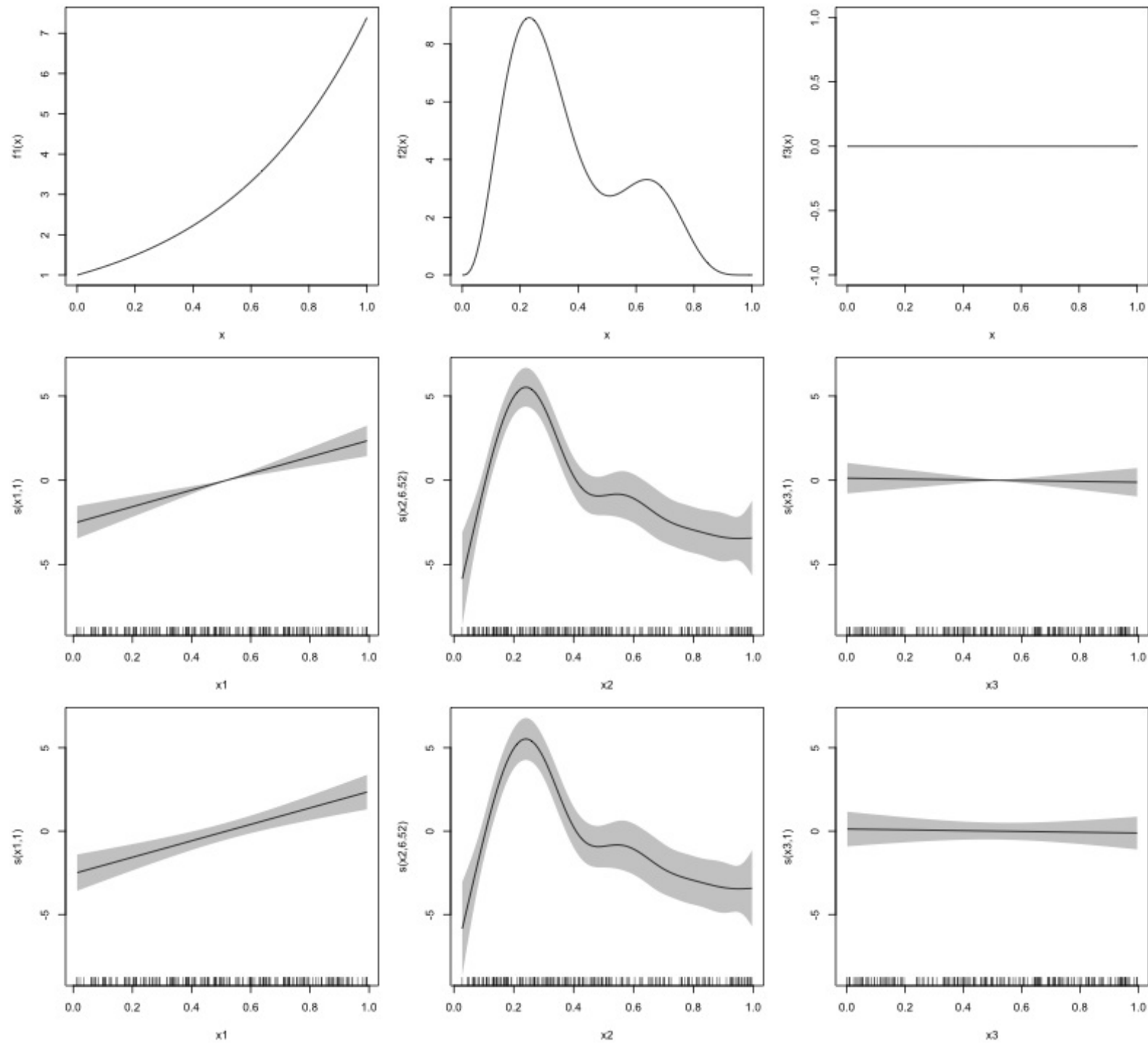
Marra & Wood (2012) suggested a solution to situation 1., namely true functions close to the penalty null space.

Smooths are normally subject to *identifiability* constraints (centred), which leads to zero variance where the estimated function crosses the zero line.

Instead, compute intervals for j th smooth as if it alone had the intercept; identifiability constraints go on the other smooth terms.

Use `seWithMean = TRUE` in call to `plot.gam()`

Example



p values for smooths

p values for smooths

...are approximate:

1. they don't really account for the estimation of λ_j — treated as known
2. rely on asymptotic behaviour — they tend towards being right as sample size tends to ∞

Also, p values in `summary.gam()` have changed a lot over time — all options except current default are deprecated as of v1.18-13.

The approach described in Wood (2006) is “*no longer recommended*”!

p values for smooths

...are a test of **zero-effect** of a smooth term

Default p values rely on theory of Nychka (1988) and Marra & Wood (2012) for confidence interval coverage.

If the Bayesian CI have good across-the-function properties, Wood (2013a) showed that the p values have

- almost the correct null distribution
- reasonable power

Test statistic is a form of χ^2 statistic, but with complicated degrees of freedom.

p values for unpenalized smooths

The results of Nychka (1988) and Marra & Wood (2012) break down if smooth terms are unpenalized.

This include i.i.d. Gaussian random effects, (e.g. bs = "re".)

Wood (2013b) proposed instead a test based on a likelihood ratio statistic:

- the reference distribution used is appropriate for testing a H_0 on the boundary of the allowed parameter space...
- ...in other words, it corrects for a H_0 that a variance term is zero.

p values for smooths

have the best behaviour when smoothness selection is done using **ML**, then **REML**.

Neither of these are the default, so remember to use `method = "ML"` or `method = "REML"` as appropriate

p values for parametric terms

...are based on Wald statistics using the Bayesian covariance matrix for the coefficients.

This is the “right thing to do” when there are random effects terms present and doesn't really affect performance if there aren't.

Hence in most instances you won't need to change the default `freq = FALSE` in `summary.gam()`

anova()

anova()

mgcv provides an `anova()` method for "gam" objects:

1. Single model form: `anova(m1)`
2. Multi model form: `anova(m1, m2, m3)`

anova() --- single model form

This differs from anova() methods for "lm" or "glm" objects:

- the tests are Wald-like tests as described for `summary.gam()` of a H_0 of zero-effect of a smooth term
- these are not *sequential* tests!

anova()

```
b1 <- gam(y ~ x0 + s(x1) + s(x2) + s(x3), method = "REML")
anova(b1)
```

Family: gaussian
Link function: identity

Formula:
y ~ x0 + s(x1) + s(x2) + s(x3)

Parametric Terms:

	df	F	p-value
x0	3	26.94	1.57e-14

Approximate significance of smooth terms:

	edf	Ref.df	F	p-value
s(x1)	1.000	1.001	26.677	5.83e-07
s(x2)	6.694	7.807	18.755	< 2e-16
s(x3)	1.000	1.000	0.068	0.795

anova() --- multi model form

The multi-model form should really be used with care — the p values are really *approximate*

```
b1 <- gam(y ~ s(x0) + s(x1) + s(x2) + s(x3) + s(x4) + s(x5), data = dat,  
          family=poisson, method = "ML")  
b2 <- update(b1, . ~ . - s(x3) - s(x4) - s(x5))  
anova(b2, b1, test = "LRT")
```

Analysis of Deviance Table

```
Model 1: y ~ s(x0) + s(x1) + s(x2)  
Model 2: y ~ s(x0) + s(x1) + s(x2) + s(x3) + s(x4) + s(x5)  
  Resid. Df Resid. Dev      Df Deviance Pr(>Chi)  
1      186.23      248.97  
2      183.34      248.01  2.8959   0.96184    0.795
```

For *general smooths* deviance is replaced by $-2 \log(\hat{\beta})$

AIC for GAMs

AIC for GAMs

- Comparison of GAMs by a form of AIC is an alternative frequentist approach to model selection
- Rather than using the marginal likelihood, the likelihood of the β_j *conditional* upon λ_j is used, with the EDF replacing k , the number of model parameters
- This *conditional* AIC tends to select complex models, especially those with random effects, as the EDF ignores that λ_j are estimated
- Wood et al (2015) suggests a correction that accounts for uncertainty in λ_j

$$\text{AIC} = -2l(\hat{\beta}) + 2\text{tr}(V_{\beta}')$$

AIC

In this example, x_3 , x_4 , and x_5 have no effects on y

```
AIC(b1, b2)
```

	df	AIC
b1	15.03493	847.7961
b2	12.12435	842.9368

References

- Marra & Wood (2011) *Computational Statistics and Data Analysis* **55** 2372–2387.
- Marra & Wood (2012) *Scandinavian journal of statistics, theory and applications* **39**(1), 53–74.
- Nychka (1988) *Journal of the American Statistical Association* **83**(404) 1134–1143.
- Wood (2006) *Generalized Additive Models: An Introduction with R*. Chapman and Hall/CRC.
- Wood (2013a) *Biometrika* **100**(1) 221–228.
- Wood (2013b) *Biometrika* **100**(4) 1005–1010.