#### **REML** estimation of penalized GLMs

**Simon Wood** Mathematical Sciences, University of Bath, U.K.

## A model for octane rating



- Aim to predict octane rating from NIR spectrum,  $g(\nu)$ .
- A possible model is that octane<sub>i</sub> ~ Gamma,

$$\mathbb{E}(\text{octane}_i) = \int f(\nu) \boldsymbol{g}_i(\nu) \boldsymbol{d}\nu$$

where *f* is a function to be estimated.

## **Diabetic retinopathy**



▶ Possible model...ret ~ Bernoulli,

$$\begin{aligned} \text{logit}\{\mathbb{E}(\texttt{ret})\} &= f_1(\texttt{dur}) + f_2(\texttt{bmi}) + f_3(\texttt{gly}) \\ &+ f_4(\texttt{dur},\texttt{bmi}) + f_5(\texttt{dur},\texttt{gly}) + f_6(\texttt{gly},\texttt{bmi}) \end{aligned}$$

#### Penalized GLMs

Examples are special cases of the model

$$g(\mu_i) = \mathbf{A}_i \mathbf{\theta} + \sum_j L_{ij} f_j, \quad y_i \sim \mathrm{EF}(\mu_i, \phi)$$

- y<sub>i</sub> is an exponential family distributed univariate response,
   A is a known parametric model matrix, with corresponding coefficients θ, some of which may be random, f<sub>j</sub> is a smooth function of one or more predictor variables, and L<sub>ij</sub> is a known linear functional.
- Such models include GAMs, GAMMs, varying coefficient models and signal regression models.

#### Practical representation

To make these models practical, it helps to represent each f<sub>i</sub> using a linear basis expansion,

$$f_j(x) = \sum_k^K \gamma_k b_k(x)$$

where the coefficients,  $\gamma_k$ , are unknown, the basis functions  $b_k(x)$  are chosen for convenience, and *K* is set large enough to avoid underfitting.

Then the model becomes

$$g(\mu_i) = \mathbf{X}_i \boldsymbol{\beta}, \quad \mathbf{y}_i \sim \mathrm{EF}(\mu_i, \phi)$$

where **X** is a model matrix containing **A** and evaluated versions of  $L_{ij}b_k$ , and  $\beta$  contains  $\theta$  and the  $f_j$  coefs.

## Estimation

► To avoid overfit, estimation is by penalized MLE.

minimise 
$$D(\beta) + \sum \lambda_k \beta^T \mathbf{S}_k \beta$$
 w.r.t.  $\beta$ 

*D* is model deviance, and the  $\lambda_k \beta^T \mathbf{S}_k \beta$  penalize overfit.

- Typically β<sup>T</sup>S<sub>k</sub>β measures 'wiggliness' of f<sub>j</sub>, or are random effect shrinkage factors.
- Smoothing parameters \u03c6<sub>k</sub> control smoothness of the f<sub>j</sub>, or random effect variances.
- There is a Bayesian/ mixed model interpretation: imposition of the penalties is equivalent to using a prior

$$\boldsymbol{\beta} \sim \boldsymbol{N}\left\{ \mathbf{0}, \left(\sum \lambda_k \mathbf{S}_k\right)^{-} \right\}$$

• Given  $\lambda$ ,  $\hat{\beta}$  is found by penalized IRLS. What about  $\lambda$ ...

#### How to estimate $\lambda$ : cross validation



- Leave out one datum at a time, fit model to rest, and find squared error in predicting the left out one.
- Average these squared errors over all data. Find λ to minimize this MSE.
- Invariant version is GCV.

## How to estimate $\lambda$ : marginal likelihood



- Draw β from prior implied by λ. Find average value of likelihood for these draws.
- Choose  $\lambda$  to maximize this average likelihood.
- Formally, maximize  $\int f(\mathbf{y}, \boldsymbol{\beta}) d\boldsymbol{\beta}$  w.r.t.  $\lambda$ .

## GCV vs. REML comparison



- Asymptotic MSE better for GCV, but GCV λ̂ converges much more slowly than REML version.
- GCV has greater tendency to multiple minima and severe undersmoothing.

### The upshot

- REML or ML based smoothness selection is often preferable to GCV or AIC in practice.
- But the computational methods for REML/ML are substantially less robust than the GCV/AIC equivalents.
- This requires some explanation...

# Computing $\hat{\boldsymbol{\beta}}|\lambda$ and then $\boldsymbol{\lambda}$

- Given  $\lambda$ , iterate a PIRLS scheme to convergence...
  - 1. Form  $z_i = \mathbf{X}_i \hat{\boldsymbol{\beta}} + (y_i \hat{\mu}_i) g'(\hat{\mu}_i) / \alpha_i$  and  $w_i = \alpha_i / \{ V(\hat{\mu}_i) g'(\hat{\mu}_i)^2 \}.$
  - 2. Minimize  $\sum w_i(z_i \mathbf{X}_i \beta)^2 + \sum_k \lambda_k \beta^T \mathbf{S}_k \beta$  w.r.t.  $\beta$  to get new estimate of  $\hat{\beta}$ .

 $\alpha_i = 1 + (y_i - \mu_i)(V'/V_i + g''_i/g'_i)$  for Newton or  $\alpha_i = 1$  for Fisher scoring.

- Denote the resulting penalized MLE by  $\hat{\beta}_{\lambda}$ .
- Two options for  $\lambda$  estimate computation...
  - 1. Single iteration: insert a GCV or REML based  $\lambda$  estimation step into the PIRLS.
  - 2. Nested iteration: let REML or GCV score depend on  $\beta$  only through  $\hat{\beta}_{\lambda}$ , so each iteratively proposed  $\lambda$  value requires a full PIRLS to find  $\hat{\beta}_{\lambda}$ .

## $\lambda$ computations compared

- Single iteration (which includes PQL) need not converge to a fixed β<sub>λ</sub>.
- ► Nested iteration is much more tedious to implement, but converges to a fixed Â<sub>λ</sub>.
- Nested iteration was only available for GCV etc, not REML etc.
- Aim is to rectify this, to make REML based penalized GLM estimation as routine and reliable as parametric GLM estimation.

## Laplace approximate REML

- Let  $f(\mathbf{y}, \beta)$  be the joint density of the data  $\mathbf{y}$  and  $\beta$ .
- To get restricted likelihood, integrate out β approximately...
- Replace f(y, β) by exponential of truncated Taylor expansion of log f, about β<sub>λ</sub>.
- Approximation is close enough to a Normal density to be integratable. Hooray!
- Twice log approximate restricted likelihood is

$$2I_r = 2I(\hat{\beta}) + \log|\mathbf{S}/\phi|_+ - \hat{\beta}^{\mathrm{T}}\mathbf{S}\hat{\beta}/\phi - \log|\mathbf{H} + \mathbf{S}/\phi| + M_{\rho}\log(2\pi).$$

**H** is Hessian of  $-\log f$ , and subscript dropped on  $\hat{\beta}_{\lambda}$ .

# Computing $\hat{\lambda}$

- ►  $I_r$  maximised w.r.t. log  $\lambda$  by Newton's method.
- Each trial log  $\lambda$  will require a full PIRLS to get  $\hat{\beta}_{\lambda}$ .
- Some computational considerations...
  - 1. log  $|\mathbf{S}/\phi|_+$ , where  $\mathbf{S} = \sum \lambda_k \mathbf{S}_k$ , is very awkward.
  - 2.  $\mathbf{H} = \mathbf{X}^{\mathrm{T}}\mathbf{W}\mathbf{X}/\phi$  if Newton PIRLS is used ( $\mathbf{W} = \operatorname{diag}(w_i)$ ).
  - Need derivatives of β<sub>λ</sub> w.r.t. log λ, to get derivatives of *I<sub>r</sub>* for Newton steps. Implicit differentiation gives these easily, if Newton based PIRLS is used.
  - 4. **But** negative weights can occur for Newton PIRLS, making stable computation difficult.

## The problem with log $|\mathbf{S}|_+$

- > Naive  $|\mathbf{S}|_+$  evaluation can go badly wrong.
- Consider log  $|\mathbf{S}_1 + \lambda \mathbf{S}_2|_+$  when  $\lambda \to \mathbf{0}$

```
> S <- S1 + S2*1e-18; S
     [,1] [,2] [,3] [,4] [,5] [,6]
[1,] 1 -1 0 0e+00 0e+00 0e+00
[2,] -1 2 -1 0e+00 0e+00 0e+00
[3,] 0 -1 1 0e+00 0e+00 0e+00

      [4,]
      0
      0
      1e-18
      -1e-18
      0e+00

      [5,]
      0
      0
      -1e-18
      2e-18
      -1e-18

[6,] 0 0 0 0e+00 -1e-18 1e-18
> sum(log(eigen(S)$values[1:4])) ## naive
[1] -73.39584
> sum(log(eigen(S1)$values[1:2])) + ## true
> sum(log(eigen(S2)$values[1:2]*1e-18))
[1] -80.69584
> eigen(S) $values ## why?
[1] 3.000000e+00 1.000000e+00 2.220446e-15 2.000000e-18
[5] 1.000000e-18 1.000000e-18
```

### Further problems with log |S|

Suppose that the non-zero sub matrices of S<sub>i</sub> need not be full rank, but overlap...

```
> S1;S2
    [,1] [,2] [,3] [,4] [,5]
                                  [,1] [,2] [,3] [,4] [,5]
[1,]
   1 -1 0
                 0
                        0
                              [1,]
                                    0
                                         0
                                             0
                                                      0
[2,] -1 2 -1
                             [2,]
                                    0
                   0
                        0
                                         0
                                             0
                                                 0
                                                      0
[3,] 0 -1 1 0 0
                             [3,] 0 0
[4,] 0 0
                                             1 0
                                                      0
[4,] 0 0 0
                   0 0
                                             0 1
                                                      0
[5,] 0 0 0
                   0
                        0
                              [5,]
                                    0
                                         Ο
                                             0
                                                 0
                                                      1
> S <- S1 + S2*1e-18
> sum(log(eigen(S)$values))
[11 - 115.5355]
> sum(log(abs(diag(qr.R(qr(S))))))
[1] -118.3859
```

 Are either of these right? No. S<sub>1</sub> is only rank 2, but its 'numerical footprint' extends beyond the first 2 columns of S, obliterating part of the rank 3 matrix, S<sub>2</sub>.

# Solution for log |S|

- Similarity transform to confine S<sub>1</sub> to a block of size corresponding to its rank.
- Consider eigen-decomposition UDU<sup>T</sup> = S<sub>1</sub>. Then |S<sub>1</sub> + λS<sub>2</sub>| = |D + λU<sup>T</sup>S<sub>2</sub>U|. If zero eigenvalues in D are set to exactly 0, then r.h.s. evaluates correctly.

```
> S
       [,1] [,2] [,3] [,4] [,5]
[1,] 1 -1 0 0e+00 0e+00
[2,] -1 2 -1 0e+00 0e+00
[3,] 0 -1 1 0e+00 0e+00
[4,] 0 0 0 1e-18 0e+00
[5,] 0 0 0 0e+00 1e-18
> es <- eigen(S1); U <- es$vectors
> D <- es$values; D[3:5] <- 0
> Sp <- diag(D) + t(U) %*%S2%*%U * lambda
> sum(log(abs(diag(qr.R(qr(Sp))))))
[1] -124.3396
```

#### ... why did that work, and does it generalize?

The large elements for S<sub>1</sub> got confined to a 2 × 2 block, in the similarity transformed version of S

> Sp

	[,1]	[,2]	[,3]	[,4]	[,5]
[1,]	3.000000e+00	-2.886751e-19	-2.357023e-19	0e+00	0e+00
[2,]	-2.886751e-19	1.000000e+00	4.082483e-19	0e+00	0e+00
[3,]	-2.357023e-19	4.082483e-19	3.333333e-19	0e+00	0e+00
[4,]	0.000000e+00	0.000000e+00	0.000000e+00	1e-18	0e+00
[5,]	0.000000e+00	0.000000e+00	0.000000e+00	0e+00	1e-18

- QR decomposition to get the determinant doesn't mix columns, so the determinant calculation is now stable.
- This approach can be generalized to more than two S<sub>j</sub> matrices (and to rank deficient S).
- The reparameterization used here also yields the most stable computation of β̂ and |X<sup>T</sup>WX + S|, so is always used (recomputed for each trial λ).

#### Newton penalized iteratively reweighted least squares

Newton PIRLS involves minimising

$$S = \sum w_i (z_i - \mathbf{X}_i \beta)^2 + \beta^{\mathrm{T}} \mathbf{S} \beta$$

where some  $w_i$  can be negative.

- Formal solution (X<sup>T</sup>WX + S)β̂ = X<sup>T</sup>Wz, where W = diag(w<sub>i</sub>), is far too badly-conditioned to use here.
- ► For positive *w<sub>i</sub>*, use QR decomposition method on

minimise 
$$S = \left\| \sqrt{\mathbf{W}} \left( \begin{bmatrix} \mathbf{z} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{X} \\ \sqrt{\mathbf{S}} \end{bmatrix} \beta \right) \right\|^2$$
 w.r.t.  $\beta$ 

If some w<sub>i</sub> < 0 use QR on</p>

minimise 
$$S = \left\| \sqrt{absW} \left( \begin{bmatrix} \mathbf{z} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{X} \\ \sqrt{\mathbf{S}} \end{bmatrix} \beta \right) \right\|^2$$
 + correction

# Derivatives of $\hat{eta}$ w.r.t. $ho = \log \lambda$

• Consider single  $\lambda$ .  $\hat{\beta}_{\lambda}$  is the solution to

$$\left. \frac{\mathrm{d}\boldsymbol{D}}{\mathrm{d}\boldsymbol{\beta}} \right|_{\hat{\boldsymbol{\beta}}_{\lambda}} + 2\lambda \mathbf{S}\hat{\boldsymbol{\beta}} = \mathbf{0}$$

• Differentiating w.r.t.  $\rho = \log \lambda$  gives

$$\frac{\mathrm{d}^2 D}{\mathrm{d}\beta \mathrm{d}\beta^{\mathrm{T}}}\bigg|_{\hat{\boldsymbol{\beta}}_{\lambda}}\frac{\mathrm{d}\hat{\boldsymbol{\beta}}_{\lambda}}{\mathrm{d}\rho}+2\lambda\mathbf{S}\hat{\boldsymbol{\beta}}+2\lambda\mathbf{S}\frac{\mathrm{d}\hat{\boldsymbol{\beta}}}{\mathrm{d}\rho}=\mathbf{0}$$

Re-arranging yields...

$$\frac{\mathrm{d}\hat{\boldsymbol{\beta}}}{\mathrm{d}\boldsymbol{\rho}} = -2\lambda \left(\frac{\mathrm{d}^2\boldsymbol{D}}{\mathrm{d}\boldsymbol{\beta}\mathrm{d}\boldsymbol{\beta}^{\mathrm{T}}} + 2\lambda \mathbf{S}\right)^{-1} \mathbf{S}\hat{\boldsymbol{\beta}}$$

... which is computable from left-overs!

# Key points in summary

- Find ρ = log(λ) by Newton's method optimization of the log restricted likelihood, *l<sub>r</sub>*.
- For each trial  $\rho$  vector ...
  - 1. Re-parameterize so that  $\log |\mathbf{S}|_+$  is computable.
  - 2. Solve for  $\hat{\beta}_{\lambda}$  by stable Newton PIRLS.
  - 3. Get the derivatives of  $\hat{\beta}_{\lambda}$  by implicit differentiation.
  - 4. Slog out the derivatives of  $l_r$  itself, from which an update for  $\hat{\rho}$  can be computed.
- Full details are exceptionally tedious, but these are the main points.

## Implementation: mgcv

- R package mgcv implements this REML based estimation method.
- Consider example GAMM.  $y_i \sim \text{Poi}$ ,

$$\log \mathbb{E}(\mathbf{y}_i) = f_1(\mathbf{x}_i)t_i + f_2(\mathbf{v}_i, \mathbf{w}_i) + \gamma_k, \quad \gamma_k \sim N(0, \sigma_g^2)$$

if observation *i* from level *k* of factor variable, group.

gam(y~s(x,by=t)+s(v,w)+s(group,bs="re"), family=poisson,method="REML")

#### Example: adaptive smoothing



- Top is normal spline smooth.
- Lower is an adaptive spline smooth, where penalty is split up, so that degree of smoothness can vary with time.

## Example: A model for octane rating



- Aim to predict octane rating from NIR spectrum,  $g(\nu)$ .
- A possible model is that octane; ~ Gamma,

$$\log \mathbb{E}(\text{octane}_i) = \int f(\nu) g_i(\nu) d\nu$$

where *f* is a function to be estimated.

#### Octane model estimates



 If nm is a matrix each row of which contains the spectral wavelengths and NIR is a corresponding matrix of spectral measurements, then model is fit by

```
gam(octane<sup>s</sup>(nm, by=NIR, bs="ad"),
family=Gamma(log))
```

## Example: Diabetic retinopathy



Possible model...ret ~ Bernoulli,

$$\begin{aligned} \text{logit}\{\mathbb{E}(\texttt{ret})\} &= f_1(\texttt{dur}) + f_2(\texttt{bmi}) + f_3(\texttt{gly}) \\ &+ f_4(\texttt{dur},\texttt{bmi}) + f_5(\texttt{dur},\texttt{gly}) + f_6(\texttt{gly},\texttt{bmi}) \end{aligned}$$

## **Retinopathy estimates**



#### Retinopathy estimates: interaction



adigreen are +/- TRUE s.e.

#### **Advert**

 Wood, S.N. (2011) Fast stable REML and ML estimation of semiparametric GLMs. JRSSB 73(1):1-34