

## The next weeks

---

March 9	§10.6 Overdispersion and quasi-likelihood, GEEs
March 16	§10.7 Semiparametric models
March 23	Generalized additive models and lasso
March 30	Finishing pieces, + review

---

Homework 3: due April 2, 5 pm

Final Test: April 17, 1 - 3 pm – posted by Monday, March 19

↑  
marks due 20th -

## Estimating functions and quasi-likelihood

- ▶ suppose we assume only that  $E(Y_j) = \mu_j(\beta)$ ,  $\text{Var}(Y_j) = \phi a_j V(\mu_j)$ , as in most glm's
- ▶ and we use the glm estimates of  $\beta$ , defined by the score equation

$$\sum_{j=1}^n \frac{y_j - \mu_j}{a_j V(\mu_j)} \frac{x_{jr}}{g'(\mu_j)} = 0 \quad (*) \quad ||$$

- ▶ n.b. Davison calls LHS  $g(Y; \beta)$ , different  $g$

- ▶ using only (\*), we have  $g(Y; \tilde{\beta}) = 0$

$$E\{g(Y; \beta)\} = 0; \quad E\left\{-\frac{\partial g(Y; \beta)}{\partial \beta}\right\} = \text{Var}\{g(Y; \beta)\}$$

- ▶ thus  $g$  has two properties in common with the the score function from a log-likelihood

## ... estimating functions and quasi-likelihood



$$g(Y; \beta) = \sum_{j=1}^n \frac{Y_j - \mu_j}{a_j V(\mu_j)} g'(\mu_j) = 0$$

$E(\cdot) = 0$

▶  $g(Y; \beta)$  is the  $\beta$ -derivative of

$$Q(\beta; Y) = \sum_{j=1}^n \int_{Y_j}^{\mu_j} \frac{Y_j - u}{\phi a_j V(u)} du,$$

▶  $w_j = 1 / \{g'(\mu_j)^2 \phi a_j V(\mu_j)\}$  as usual, have assumed

▶ Quasi-likelihood estimator  $\tilde{\beta} \sim N(\beta, (X^T W X)^{-1})$  where  $W = \text{diag}(w_1, \dots, w_n)$

▶ inflate estimated standard errors for  $\tilde{\beta}_j$  by  $\hat{\phi}^{1/2}$



relative bias.

$$\hat{\phi} = \frac{1}{n-p} \sum_{j=1}^n \frac{(y_j - \hat{\mu}_j)^2}{a_j V(\hat{\mu}_j)}$$

## ... estimating functions and quasi-likelihood

```
> toxo.glm1 = glm(cbind(r,m-r) ~ rain + I(rain^2) + I(rain^3),  
family = quasibinomial)  
> summary(toxo.glm1)
```

...

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-2.902e+02	1.215e+02	-2.388	0.0234	*
rain	4.500e-01	1.876e-01	2.398	0.0229	*
I(rain^2)	-2.311e-04	9.616e-05	-2.404	0.0226	*
I(rain^3)	3.932e-08	1.635e-08	2.405	0.0225	*

Null dev.      resid dev.       $\hat{\phi} = 1.94$

```
> (74.212 - 62.635)/3/1.94
```

```
[1] 1.989175
```

```
> pf(.Last.value, 3, 30, lower.tail = F)
```

```
[1] 0.1368155
```

$\tilde{\beta} \rightarrow N(\beta; \cdot)$

## ... estimating functions and quasi-likelihood

- ▶ even if  $V(\mu)$  is incorrectly specified,  $\tilde{\beta}$  is still consistent

$$E\left(-\frac{\partial g}{\partial \beta}\right)$$

$$E\left(-\frac{\partial g}{\partial \beta}\right)$$

▶

$$a. \text{Var}(\tilde{\beta}) = (X^T W X)^{-1} \text{Var}\{g(Y; \beta)\} (X^T W X)^{-1}$$

- ▶ often is well approximated by  $(X^T W X)^{-1}$  in any case

- ▶ when extended to dependent data, called **generalized estimating equation** method

- ▶ reference: Liang & Zeger (1986, Biometrika)

← if  
right  
= I

sandwich est. of variance

## Example: estimation of spatial intensity (Yongtao Guan, Mar. 15)

- ▶ Poisson process on spatial area  $W$ , indexed by spatial locations  $s$   $-(\text{long, lat})$
- ▶  $N(s)$  counts "events" at location  $s$ ,  $N(s) \sim \text{Pois}\{\lambda(s)\}$   
 $(s, s+ds)$   $=$
- ▶ generalized linear model  $\lambda(s) = \exp\{Z(s)^T \beta\}$   $\pi$  species elev soil
- ▶ introduce correlation by assuming two points  $s_1$  and  $s_2$  have a joint intensity function
- ▶  $\lambda_2(s_1, s_2) = \lambda(s_1)\lambda(s_2)g(\|s_1 - s_2\|)$
- ▶ estimation using mean  $\lambda(\cdot)$  and variance  $\lambda_2(\cdot, \cdot)$  only as in GEE

# Generalized estimating equations $\phi V(\mu_j; \alpha_j$

- $n$  ind't  $Y_j$ 's
- ▶  $Y_j = (Y_{j1}, \dots, Y_{jn_j})$ ;  $E(Y_j) = \mu_j$ ;  $\text{var}(Y_j) = \cancel{V} V(\mu_j; \alpha)$   
↖  $n_j \times n_j$  matrix
- ▶ estimating equation for  $\beta$ :

(canonical link) ←

$$\sum_{j=1}^n \left( \frac{\partial \mu_j}{\partial \beta^T} \right) V(\mu_j; \alpha)^{-1} (Y_j - \mu_j) = \underline{0}_{p \times 1}$$

$X^T$

GLM type eg -

- ▶ multivariate version of quasi-likelihood equation ←
- ▶ needs some specification of  $V(\cdot; \cdot)$  called “working covariance matrix”
- ▶ `gee` in `library(gee)` offers several choices: independent, exchangeable,  $AR(p)$ , etc.
- \* ▶ estimate of  $\beta$  is consistent, even if  $V(\cdot; \cdot)$  is mis-specified
- ▶ but estimates of  $\text{Var}(\tilde{\beta})$  will be incorrect if  $V$  is “
- ▶ there is no quasi-likelihood that corresponds to this more general model

## Dependence through random effects

- ▶ Example: longitudinal data

$$C_1 \quad 9 \quad \forall_j \in \mathbb{R}$$

- ▶  $Y_j = (Y_{j1}, \dots, Y_{jn_j})$  vector of observations on  $j$ th individual
- ▶ recall random effects model (normal theory):

$$Y_j = X_j \beta + Z_j b_j + \epsilon_j; \quad b_j \sim N(0, \sigma^2 \Omega_b), \epsilon_j \sim N(0, \sigma^2 \Omega_j)$$

- ▶ marginal distribution:

$$Y_j \sim N(X_j \beta, \sigma^2 \Upsilon_j^{-1})$$

- ▶ sample of  $n$  i.i.d. such vectors leads to

$$Y \sim N(X \beta, \sigma^2 \Upsilon^{-1}), \quad \Upsilon^{-1} = (\Omega + Z \tilde{\Omega}_b Z^T)$$

- ▶  $\Omega = \text{diag}(\Omega_1, \dots, \Omega_n), \quad \tilde{\Omega}_b = \text{diag}(\Omega_b, \dots, \Omega_b),$

- ▶  $\Upsilon^{-1} = \text{diag}(\Upsilon_1^{-1}, \dots, \Upsilon_n^{-1})$

$$\begin{array}{l} \Upsilon \quad (n_j) \\ (Y_{j1}, \dots, Y_{jn_j}) \end{array}$$




# Generalized linear mixed models

- ▶ simplify as in last slide to canonical link  $(\theta_j = \eta_j)$



$$f(y_j | \theta_j, \phi) = \exp\left\{\frac{y_j \theta_j - b(\theta_j)}{\phi a_j} + c(y_j; \phi a_j)\right\}$$

- ▶ random effects

$$\theta_j = \underbrace{x_j^T}_{\text{fixed}} \underbrace{\beta}_{\text{fixed}} + \underbrace{z_j^T}_{\text{fixed}} \underbrace{b}_{\text{random}}, \quad b \sim N(0, \Omega_b)$$


- ▶ likelihood

$$L(\beta, \phi; y) = \prod_{j=1}^n \int f(y_j | \beta, b, \phi) f(b; \Omega_b) db$$

## ... generalized linear mixed models

- ▶ likelihood

$$L(\beta, \phi; y) = \prod_{j=1}^n \int f(y_j | \beta, \mathbf{b}, \phi) f(\mathbf{b}; \Omega_b) d\mathbf{b}$$

*simplify*

- ▶ doesn't simplify unless  $f(\cdot)$  is normal
- ▶ solutions proposed include
  - ▶ numerical integration, e.g. by quadrature
  - ▶ integration by MCMC
  - ▶ penalized quasi-likelihood – use Laplace approximation to the integral

$$\prod_{j,k,k'} f_2(y_{jk}, y_{jk'})$$

*(Composite)*

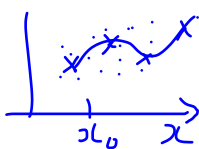
- ▶ reference: MASS library and book (§10.4):

`glmmN0`, `GLMMGibbs`, `glmmPQL`, *all in library(MASS)*  
`glmer` *in library(lme4)*

- ▶ see also Faraway (*Extending the Linear Model with R*), Ch. 10

# Semiparametric Regression §10.7

- ▶ model  $y_j = g(x_j) + \epsilon_j$ ,  $j = 1, \dots, n$   $x_j$  scalar



- ▶ mean function  $g(\cdot)$  assumed to be “smooth”

$$E(y|x) = g(x)$$

- ▶ introduce a **kernel function**  $w(u)$  and define a set of weights

$$w_j = \frac{1}{h} w\left(\frac{x_j - x_0}{h}\right) \quad 1) \quad w(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}u^2}$$

- ▶ estimate of  $g(x_0)$ :

$$\hat{g}(x_0) = \frac{\sum_{j=1}^n w_j y_j}{\sum_{j=1}^n w_j}$$

$$= \frac{e^{-\frac{1}{2}u^2}}{\sqrt{2\pi}}$$

- ▶ Nadaraya-Watson estimator (10.40) – local averaging

(tri cube)  $2) (1 - |u|^3)^3 \quad |u| \leq 1$

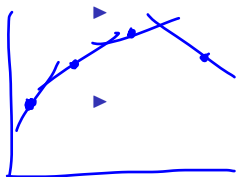
## ... kernel smoothing

- ▶ better estimates can be obtained using local regression at point  $x$

$$\frac{\sum w_i y_i}{\sum w_i}$$



$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & (x_1 - x_0) & \cdots & (x_1 - x_0)^k \\ \vdots & \vdots & & \vdots \\ 1 & (x_n - x_0) & \cdots & (x_n - x_0)^k \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}, \quad \Leftarrow$$



$$\hat{\beta} = (X^T W X)^{-1} X^T W y$$

$$W = \text{diag}(w_j)$$

$$\hat{g}(x_0) = \hat{\beta}_0$$

$$w_j = \frac{1}{h} w\left(\frac{x_0 - x_j}{h}\right)$$

- ▶ usually obtain estimates  $\hat{g}(x_j), j = 1, \dots, n$

## ... kernel smoothing

- ▶ odd-order polynomials work better than even; usually local linear fits are used
- ▶ kernel function is often a Gaussian density, or the tricube function (10.37)  $h$
- ▶ choice of bandwidth controls smoothness of function
- ▶ kernel estimators are biased
- ▶ larger bandwidth = more smoothing – increases bias, decreases variance
- ▶ some smoothers allows variable bandwidth depending on density of observations near  $x_0$
- ▶ `ksmooth` computes local averages; `loess` computes local linear regression (robustified)

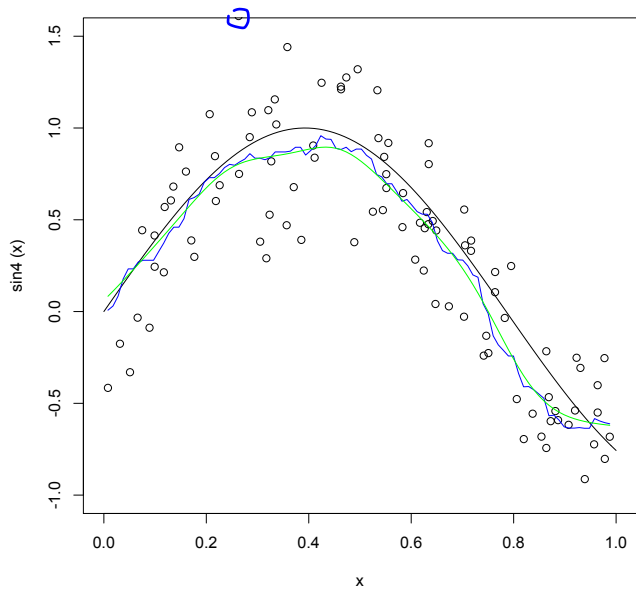
$$E \hat{g}(x_0) \neq g(x_0)$$

## Example: weighted average

```
ksmooth(x,y,kernel=c("box","normal"),bandwidth=0.5,  
        range.x=range(x),  
        n.points=max(100,length(x)), x.points)
```

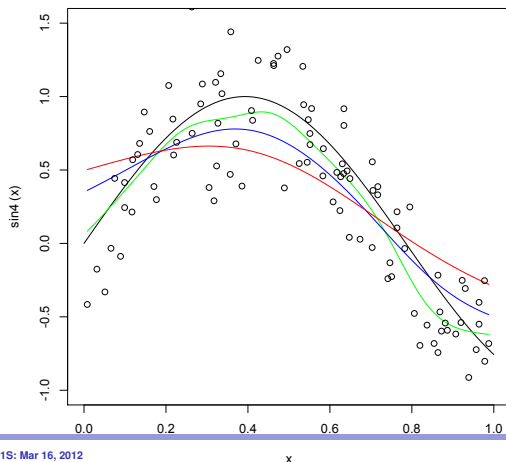
```
> eps<-rnorm(100,0,1/3)  
> x<-runif(100)  
> sin4 <- function(x){sin(4*x)}  
> y<-sin4(x)+eps  
> plot(sin4,0,1,type="l",ylim=c(-1.0,1.5),xlim=c(0,1))  
> points(x,y)  
> lines(ksmooth(x,y,"box",bandwidth=.2),col="blue")  
> lines(ksmooth(x,y,"normal",bandwidth=.2),col="green")
```





## ... Example

```
> plot(sin4, 0, 1, type="l", ylim=c(-1.0, 1.5), xlim=c(0, 1))  
> lines(ksmooth(x, y, "normal", bandwidth=.2), col="green")  
> lines(ksmooth(x, y, "normal", bandwidth=0.4), col="blue")  
> lines(ksmooth(x, y, "normal", bandwidth=0.6), col="red")
```





## Fitting in R

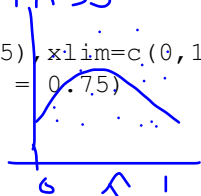
- ▶ `scatter.smooth` fits a loess curve to a scatter plot
- ▶ `loess` takes a family argument: `family = gaussian` gives weighted least squares using  $K_\lambda$  as weights and `family=symmetric` gives a robust version using Tukey's biweight
- ▶ `supsmu` implements "Friedman's super smoother": a running lines smoother with elaborate adaptive choice of bandwidth
- ▶ Library `KernSmooth` has `locpoly` for local polynomial fits, and by setting `degree = 0` gives a kernel smooth
- ▶ as usual more smoothing means larger "bias", smaller "variance"

$$W(u) = e^{-\frac{1}{2}u^2/\sqrt{u}}$$

## Example: local linear smoothing

MASS

```
> plot(sin4,0,1,type="l",ylim=c(-1,1.5),xlim=c(0,1), xlab="x", ylab="y")
> lo1 = loess(y ~ x, degree = 1, span = 0.75)
```



```
> attributes(lo1)
```

```
$names
```

[1] "n"	"fitted"	"residuals"	"enp"	"s"
[7] "two.delta"	"trace.hat"	"divisor"	"pars"	"kd"
[13] "terms"	"xnames"	"x"	"y"	"we"

```
$class
```

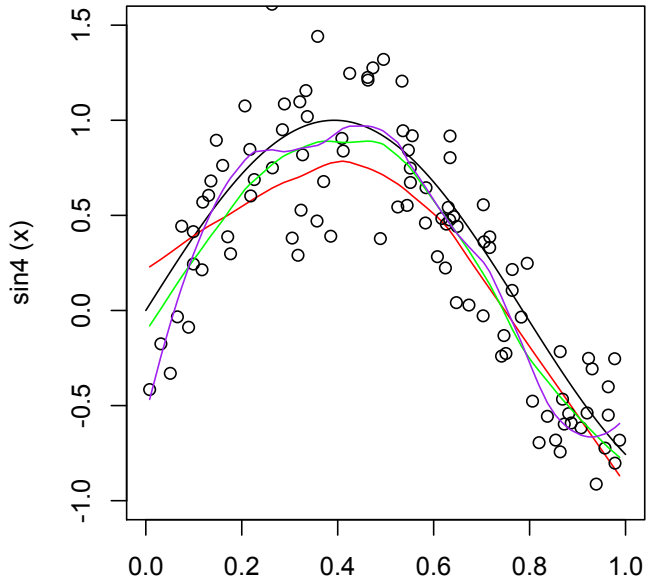
```
[1] "loess"
```

(ord = ?)

runif(0,1)

ord = order(lo1\$x)

```
> lines(lo1$x[ord], lo1$fitted[ord], col="red")
> lo2 = loess(y~x, degree=1, span=0.4)
> lo3 = loess(y~x, degree=2, span=0.4)
> lines(lo1$x[ord], lo2$fitted[ord], col="green")
> lines(lo1$x[ord], lo3$fitted[ord], col="purple")
```



## Scatter Plot with Smooth Curve Fitted by Loess

### Description

Plot and add a smooth curve computed by `loess` to a scatter plot.

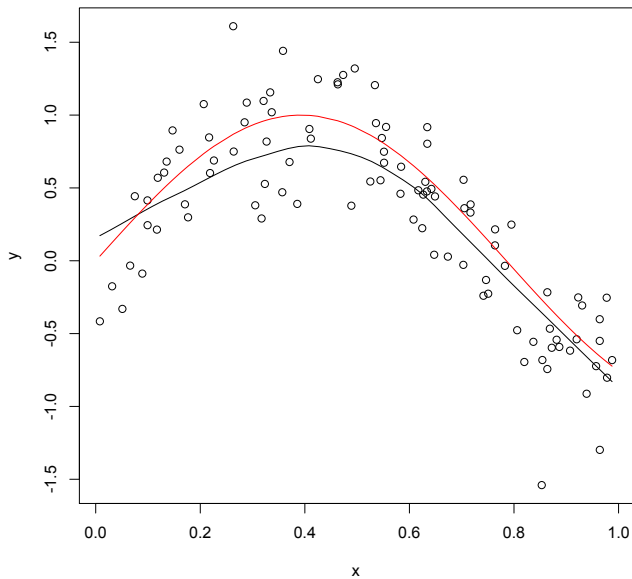
### Usage

```
scatter.smooth(x, y = NULL, span = 2/3, degree = 1,  
              family = c("symmetric", "gaussian"),  
              xlab = NULL, ylab = NULL,  
              ylim = range(y, prediction$y, na.rm = TRUE),  
              evaluation = 50, ...)
```

```
loess.smooth(x, y, span = 2/3, degree = 1,  
             family = c("symmetric", "gaussian"), evaluation = 50, ...)
```

### Arguments

- |                     |  |
|---------------------|--|
| <code>x, y</code>   | the <code>x</code> and <code>y</code> arguments provide the <code>x</code> and <code>y</code> coordinates for the plot. Any reasonable way of the coordinates is acceptable. See the function <a href="#">xy.coords</a> for details. |
| <code>span</code>   | smoothness parameter for <code>loess</code> .  |
| <code>degree</code> | degree of local polynomial used.   |
| <code>family</code> | if <code>"gaussian"</code> fitting is by least-squares, and if <code>family="symmetric"</code> a re-descending M is used.  |
| <code>xlab</code>   | label for <code>x</code> axis.   |
| <code>ylab</code>   | label for <code>y</code> axis.   |



## Friedman's SuperSmoother

### Description

Smooth the (x, y) values by Friedman's 'super smoother'.

### Usage

```
supsmu(x, y, wt, span = "cv", periodic = FALSE, bass = 0)
```

### Arguments

**x** x values for smoothing

**y** y values for smoothing

**wt** case weights, by default all equal

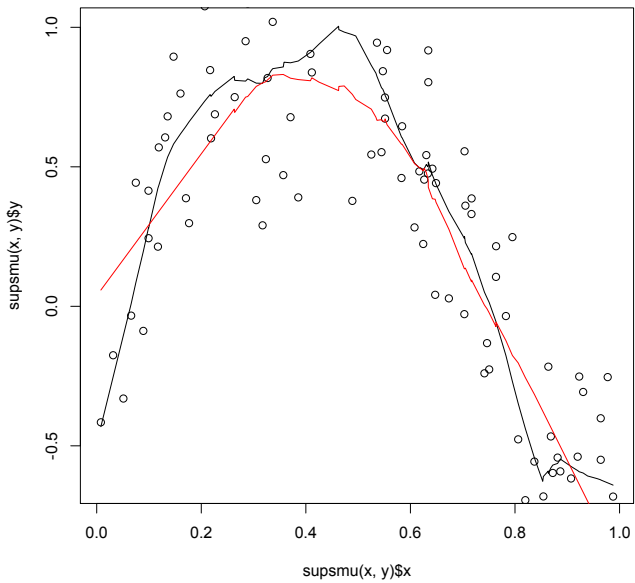
**span** the fraction of the observations in the span of the running lines smoother, or "cv" to choose the leave-one-out cross-validation.

**periodic** if TRUE, the x values are assumed to be in [0, 1] and of period 1.

**bass** controls the smoothness of the fitted curve. Values of up to 10 indicate increasing smoothness.

### Details

`supsmu` is a running lines smoother which chooses between three spans for the lines. The running lines smoothers are symmetric, with  $k/2$  data points each side of the predicted point, and values of  $k$  as  $0.5 * n$ ,  $0.2 * n$ ,  $0.05 * n$ , where  $n$  is the number of data points. If `span` is specified, a single smoother with span `span * n` is used.



# Inference from smooth functions

▶  $\hat{\beta} = (X^T W X)^{-1} X^T W y$  ( $x_0$ )

$(10.39)$   
 $n - 2\nu_1 + \nu_2$

▶  $W = \text{diag}(W_1, \dots, W_n)$

$W_j = \frac{1}{h} W\left(\frac{x_0 - x_j}{h}\right)$

▶  $\hat{g}(x_0) = \sum_{j=1}^n S(x_0; x_j, h) y_j \equiv \hat{\beta}_0$  use  $\hat{\beta}_0$  from regn

▶  $S(x_0; x_1, h), \dots, S(x_0; x_n, h)$  first row of "hat" matrix  
 $(X^T W X)^{-1} X^T W$

▶  $E\{\hat{g}(x_0)\} = \sum_{j=1}^n \underline{S(x_0; x_j, h)} g(x_j) \neq g(x_0)$

▶  $\text{var}\{\hat{g}(x_0)\} = \sigma^2 \sum_{j=1}^n S(x_0; x_j, h)^2$

$S_h = n \times n$

$\hat{\beta} = H y$

similarly  $\hat{g} = (\hat{g}(x_1), \dots, \hat{g}(x_n)) = S_h y$

$S_{ij}^h = S(x_i; x_j, h)$

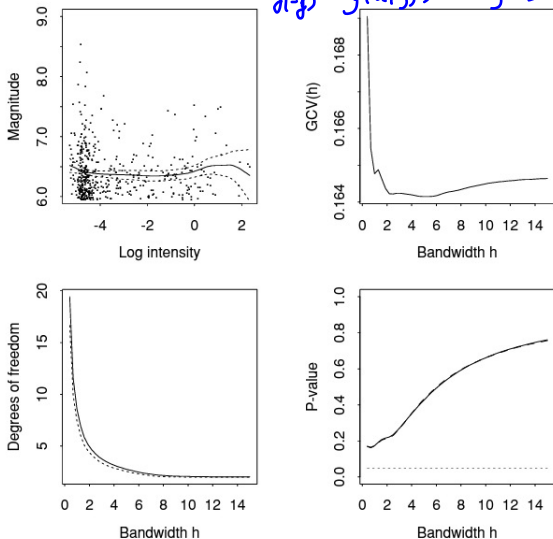
▶  $\nu_1 \approx 2\nu_1 - \nu_2$   
 $\nu_1 = \text{tr}(S_h), \nu_2 = \text{tr}(S_h^T S_h)$  suggested as

$H^T H = I$



omit  $y_j$ 

$$y_j - \hat{y}_j \sim g(x_j, z_j) \rightarrow \hat{y}_j \rightarrow \hat{y}_j \quad \sum (\hat{y}_j - y_j)^2 \leftarrow$$



**Figure 10.16** Smooth analysis of earthquake data. Upper left: local linear regression of magnitude on log intensity just before quake (solid), with 0.95 pointwise confidence bands (dots). Upper right: generalized cross-validation criterion  $GCV(h)$  as a function of bandwidth  $h$ . Lower left: relation between degrees of freedom  $v_1$  (solid),  $v_2$  (dots), and  $h$ . Lower right: significance traces for test of no relation between magnitude and log intensity, based on chi-squared approximation (dots) and saddlepoint approximation (solid). The horizontal line shows the conventional 0.05 significance level.

## Extension

- ▶ original model  $y_j = g(x_j) + \epsilon_j$

← gen LS

- ▶ extend to  $y_j \sim f(\cdot; \beta, x_j)$

gen GLM

▶

$$\max_{\beta} \sum \log f(y_j; \beta, x_j) \rightarrow \max_{\beta} \sum \frac{1}{h} w \left( \frac{x_j - x_0}{h} \right) \log f(y_j; \beta, x_j)$$

- ▶ local likelihood fitting

local

- ▶ more than 1 covariate

$$g_1(x_{1j}, x_{2j})$$

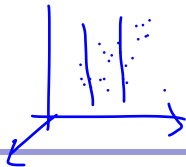
likelihood

- ▶  $E(Y_j) = g_1(x_{1j}) + g_2(x_{2j}) + \dots + g_p(x_{pj})$

- ▶ or

GAM

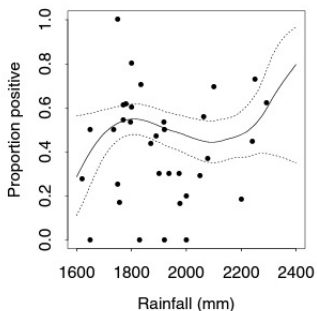
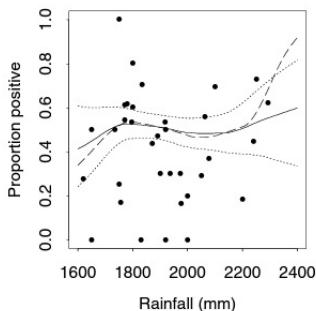
$$h(EY_j) \uparrow = \log \text{logit} \dots$$



# Example 10.32

528

10 · Nonlinear Regression Models



**Figure 10.17** Local fits to the toxoplasmosis data. The left panel shows fitted probabilities  $\hat{\pi}(x)$ , with the fit of local linear logistic model with  $h = 400$  (solid) and 0.95 pointwise confidence bands (dots). Also shown is the local linear fit with  $h = 300$  (dashes). The right panel shows the local quadratic fit with  $h = 400$  and its 0.95 confidence band. Note the increased variability due to the quadratic fit, and its stronger curvature at the boundaries.

# Flexible modelling using basis expansions

(§10.7.2)

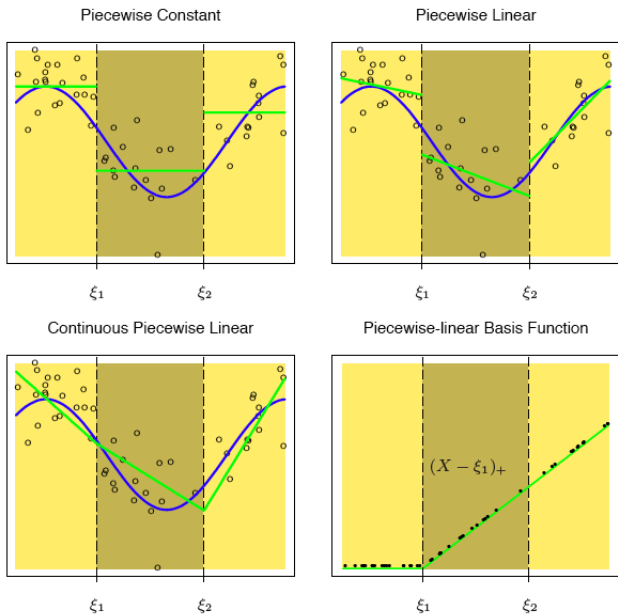
- ▶  $y_j = g(x_j) + \epsilon_j$
- ▶ Flexible linear modelling

$$g(x) = \sum_{m=1}^M \beta_m h_m(x)$$

- ▶ This is called a **linear basis expansion**, and  $h_m$  is the  $m$ th basis function
- ▶ For example if  $X$  is one-dimensional:  
 $g(x) = \beta_0 + \beta_1 x + \beta_2 x^2$ , or  
 $g(x) = \beta_0 + \beta_1 \sin(x) + \beta_2 \cos(x)$ , etc.
- ▶ Simple linear regression has  $h_1(x) = 1$ ,  $h_2(x) = x$

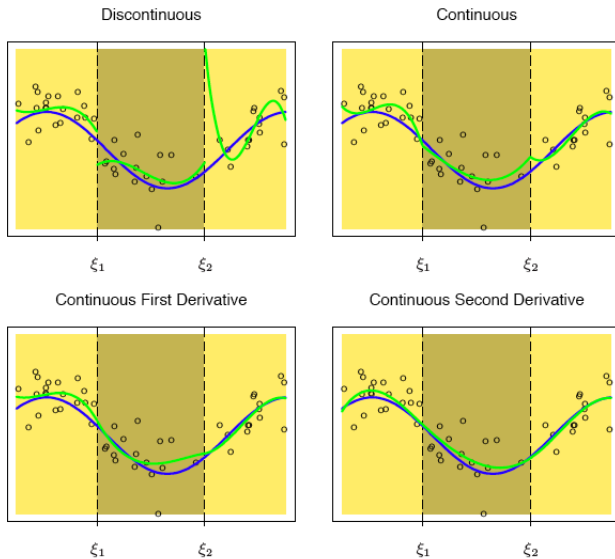
# Piecewise polynomials

- ▶ piecewise constant basis functions
$$h_1(x) = I(x < \xi_1), \quad h_2(x) = I(\xi_1 \leq x < \xi_2),$$
$$h_3(x) = I(\xi_2 \leq x)$$
- ▶ fitting by local averaging
  
- ▶ piecewise linear basis functions , with constraints
$$h_1(x) = 1, \quad h_2(x) = x$$
$$h_3(x) = (x - \xi_1)_+, \quad h_4(x) = (x - \xi_2)_+$$
- ▶ windows defined by **knots**  $\xi_1, \xi_2, \dots$
  
- ▶ **piecewise cubic basis functions**
$$h_1(x) = 1, h_2(x) = x, h_3(x) = x^2, h_4(x) = x^3$$
- ▶ continuity  $h_5(x) = (x - \xi_1)_+^3, \quad h_6(x) = (x - \xi_2)_+^3$
- ▶ continuous function, continuous first and second derivatives



**FIGURE 5.1.** The top left panel shows a piecewise constant function fit to some artificial data. The broken vertical lines indicate the positions of the two knots  $\xi_1$  and  $\xi_2$ . The blue curve represents the true function, from which the data were

## Piecewise Cubic Polynomials



**FIGURE 5.2.** A series of piecewise-cubic polynomials, with increasing orders of continuity.





