Methods of Applied Statistics I

STA2101H F LEC9101

Week 11

November 30 2022



4:56 AM · Nov 29, 2022 · Twitter for iPad

link

12 · IMS Bulletin

Volume 51 · Issue 8

Written by Witten: Impostors Anonymous



Contributing Editor Daniela Witten writes:

I was in my second year of grad school when I first heard of "impostor syndrome", the well-studied psychological phenomenon by which highly talented and accomplished people doubt their talent and

accomplishments, and live in constant fear that the outside world will discover them as frauds. I remember marveling at the possibility that some of the breathtakingly brilliant statisticians in my department might question their own abilities—had they no self-awareness?? The absurdity of their impostor syndrome stood in stark contrast to what I knew to be true: that the people around me had their figurative ducks all in a row, whereas I *uus the real fraud*.

Friends, it did not occur to me that I might, in fact, suffer from impostor syndrome until several years into my extremely hard. If all of my accomplishments to date were due to luck and/or trickery, then I had better hurry to accomplish as much as possible before my luck changed and/or my trickery was discovered! 1 am certain that I would not have achieved the same level of success as early in my career without my impostor syndrome. But, I might have been a lot happier and 90% as successful, and I firmly believe that this would have been enough, for any reasonable definition of "enough".

(I also acknowledge that impostor syndrome can manifest in different ways. For instance, some people might find themselves unable to complete a research paper due to a fear that others will discover them to be a fraud.)

Over the years, Ive learned that my impostor syndrome places a burden on those around me. If I believe that everyone else is smarter than me, then I will have unrealistically high expectations for others. This manifests not only in thinking that all of my grad students are brilliant (and in fact, they are!) but also in expecting them to constantly howe brilliant ideas which is clearly a bizer and unrealistic

Today

- 1. Upcoming events
- 2. Project
- 3. Recap
- 4. Nonparametric regression

Project due December 19 (11.59), no extensions So think of it as due on December 16 :)

Preliminary versions accepted for feedback up to Dec 11

Applied Statistics I November 30 2022

Project Guidelines	STA 2101F: Methods of Applied Statistics I 2022
Outline	
 Part I 3–5 pages, non-technical 	$12\ {\rm point}$ type, $1.5\ {\rm vertical}\ {\rm spacing},\ {\rm thank}\ {\rm you}$
 a description of the scientific problem of how (and why) the data being analyzed preliminary description of the data (plots non-technical summary for a non-statistic 	interest was collected s and tables) cian of the analysis and conclusions
 Part II 3–5 pages, technical 	LaTeX or R markdown; submit .Rmd and .pdf files
 models and analysis summary for a statistician of the analysis 	s and conclusions
 Part III Appendix R script or .Rmd file; additional plots; add 	submit .Rmd and .pdf or .html files itional analysis; References
Project Marking	
 40 points total 	
 Part I: description of data and scientific problem 5 suitability of plots and tables 5 quality of the presentation 5 	clear, non-technical, concise but thorough
 Part II: summary of the modelling and methods 5 suitability and thoroughness of the analysis 	justification for choices 10 model checks, data checks
 Part III: 	

relevance of additional material 5

complete and reproducible submission 5

Upcoming

- December 1 11.00-4.30 PostDoc Day Room 9014, Hydro Building Register here
- December 5 3.30 Data Science ARES Amy Kuceyeski, Cornell "Quantitative modeling of brain-behavior relationships" online
- December 8 3.30-4.30 Statistics Seminar Room 9014, Hydro Building Emma Hubert, Princeton "Continuous-Time Incentives in Hierarchies"









- proportional hazards regression; estimation of survivor function; interpretation of coefficients
- factorial experiments, fixed and random effects, nested and crossed factors
- · expected mean squares and components of variance
- fitting with lme4::lmer gives estimates of fixed and random effects
- general formulation of mixed models
- random effects are useful when the factor levels are not themselves meaningful, but rather a sample of possible levels
- random effects can be used to induce dependence between measurements on the same unit

Recap: Components of variance

Some factorial models:

$$y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

$$\mathbf{y}_{ijk} = \mu + \alpha_i + \beta_{ij} + \epsilon_{ijk}$$

$$\mathbf{y}_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}$$

$$y_{ij} = \beta_0 + \beta_1 x_i + b_0 + b_1 x_i + \epsilon_{ij}$$

$$y = X\beta + Zb + \epsilon$$

	meanses	cses	sector	mathach	ses	school	
	-0.43438	-0.45362	Public	1.527	-0.888	1224	9
E 8. W show constate regressions	0.71200	0.62000	Catholic	18.496	1.332	1433	248
for each school and then a more parsimonious	-0.32973	0.39173	Public	6.415	0.062	2467	1094
model with fixed effects for SES and sector	-0.13765	1.07965	Catholic	11.437	0.942	2629	1195
model with fixed effects for SES and Sector	-0.96443	-0.12357	Public	-0.763	-1.088	2639	1283
and random slope for student-level SES	-0.64918	0.36118	Public	13.156	-0.288	3657	2334
	0.40200	0.39000	Catholic	14.500	0.792	4042	2783
	0.40200	0.08000	Catholic	3.687	0.482	4042	2806
	-0.09400	1.33600	Catholic	20.375	1.242	4223	2886
	-0.10714	-0.07086	Catholic	15.550	-0.178	4511	3278
	-0.10714	0.44914	Catholic	7.447	0.342	4511	3317
	0.82498	0.07702	Catholic	18.802	0.902	5404	3656
	-0.09012	0.53212	Public	23.591	0.442	7232	5180
	0.07823	-1.17623	Public	-1.525	-1.098	7276	5223
	0.29700	-0.80500	Catholic	16.114	-0.508	7332	5278
	-0.08936	-0.08864	Catholic	be200325	-00178	atist7 36 4	Appli 5467 a
	0.15513	-0.38313	Public	18.463	-0.228	8707	6292

6

• design: one factor with I levels; J responses at each level

model

$$\mathbf{y}_{ij} = \mu + \alpha_i + \epsilon_{ij}, \quad j = 1, \dots J; \ i = 1, \dots I; \quad \epsilon_{ij} \sim (\mathbf{0}, \sigma_{\epsilon}^2)$$

Analysis of variance table

Term	degrees of freedom	sum of squares	mean square	F-statistic
treatment	(<i>I</i> − 1)	$\sum_{ij}(ar{y}_{i.}-ar{y}_{})^2$	$\sum_{ij}(ar{y}_{i.}-ar{y}_{})^2/(l-1)$	MS _{treatment} /MS _{error}
error	I(J - 1)	$\sum_{ij} (y_{ij} - \bar{y}_{i.})^2$	$\sum_{ij} (y_{ij} - \bar{y}_{i.})^2 / \{I(J-1)\}$	
total(corrected)	lJ — 1	$\sum_{ij} (y_{ij} - \bar{y}_{})^2$		

$$\begin{split} \mathrm{E}(\mathsf{SS}_{error}) &= \mathsf{I}(\mathsf{J}-\mathsf{1})\sigma_{\epsilon}^{2},\\ \mathrm{E}(\mathsf{SS}_{treatment}) &= (\mathsf{I}-\mathsf{1})(\mathsf{J}\sigma_{\alpha}^{2}+\sigma_{\epsilon}^{2}) \end{split}$$

Expected Mean Squares

$$\begin{split} \mathrm{E}(SS_{error}) &= I(J-1)\sigma_{\epsilon}^{2},\\ \mathrm{E}(SS_{treatment}) &= (I-1)(J\sigma_{\alpha}^{2}+\sigma_{\epsilon}^{2}) \end{split}$$

- more usual to have a model with some fixed effects: treatments, explanatory variables (age, income, ...)
- and some random effects: cluster, family, school, hospital, ...
- the general form of a linear mixed effect model is

 $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\epsilon}; \quad \boldsymbol{\epsilon} \sim \mathbf{N}(\mathbf{O}, \sigma^2 \mathbf{I}_n), \boldsymbol{\gamma} \sim \mathbf{N}(\mathbf{O}, \sigma^2 \mathbf{D}) \implies \mathbf{y} \sim \mathbf{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2(\mathbf{I}_n + \mathbf{Z}\mathbf{D}\mathbf{Z}^T))$

- model matrix $X_{n \times p}$, fixed effects β ; model matrix $Z_{n \times q}$, random effects γ
- fit with maximum likelihood for β , "REML" for variance components
- rely on 1me4 for inference about fixed and random effects

```
> summary(rat.mixed}
Linear mixed model fit by REML ['lmerMod']
Formula: y ~ week + (week | rat)
Data: rat.growth
```

```
• • •
```

Random effects:

	Groups	Name		Variand	ce Sto	d.Dev.	Corr
	rat	(Inte	ercept) 119.54	10	.933	
		week		12.49	3	.535	0.18
	Residual	L		33.84	5	.817	
1	Number of	obs:	150,	groups:	rat,	30	

Fixed effects:

	Estimate	Std. Error	t	value	
(Intercept)	156.0533	2.1590		72.28	
week	43.2667	0.7275		59.47	

"the estimated mean weight in week 1 is 156, but the variability from rat to rat has standard deviation of about 11 about this.

The slopes show similarly large variation.

The measurement error variance $\hat{\sigma}^2 = 5.82^2$ is smaller than the inter-rat variation in intercepts but exceed that for slopes"

Nonparametric Regression

- model $y_i = f(x_i) + \epsilon_i$, i = 1, ..., n x_i scalar
- mean function $f(\cdot)$ assumed to be "smooth"
- introduce a kernel function K(u) and define a set of weights

$$W_i = \frac{1}{\lambda} K\left(\frac{X_i - X_0}{\lambda}\right)$$

• estimate of f(x), at $x = x_0$:

$$\hat{f}_{\lambda}(\mathbf{x}_{0}) = \frac{\sum_{i=1}^{n} w_{i} y_{i}}{\sum_{i=1}^{n} w_{i}}$$

Nadaraya-Watson estimator – local averaging

local polynomial of degree o

Applied Statistics I November 30 2022

Kernel smoothers

- choice of bandwidth, λ controls smoothness of function
- larger bandwidth = more smoothing
- kernel estimators are biased
- making the estimate smoother increases bias, decreases variance
- choice of kernel function, *K*(·), controls smoothness and "local-ness"
- Faraway recommends Epanechnikov kernel $K(x) = \frac{3}{4}(1-x^2), |x| \le 1$
- ksmooth(base) offers only uniform (box) or normal
- bkde(KernSmooth) offers normal, box, epanech, biweight, triweight
- biweight: $K(x) \propto (1-|X|^2)^2, |x| \leq 1$ triweight: $K(x) \propto (1-|X|^2)^3, |x| \leq 1$

Examples



exb <- data.frame(exb)</pre>

```
plota <- ggplot(exa) + geom_point(aes(x,y)) +
geom_line(aes(x,m))+ ggtitle("Example A")</pre>
```

```
plotb <- ggplot(exb) + geom_point(aes(x,y)) +
geom_line(aes(x,m))+ ggtitle("Example B")</pre>
```

```
plotc <- ggplot(faithful) + geom_point(aes(eruptions,waiting)) +
ggtitle("Old Faithful")</pre>
```

grid.arrange(plota, plotb, plotc, nrow=1) #in gridExtra library

... Examples



Applied Statistics I November 30 2022

```
with(faithful, plot(eruptions, waiting, cex=2, main = "bandwidth=0.1", pch="."))
lines(locpoly(faithful$eruptions,faithful$waiting,drv=0L,
    degree=0,bandwidth=.1), col = "blue")
```

with(faithful, plot(eruptions, waiting, cex=2, main = "bandwidth=0.5", pch="."))
lines(locpoly(faithful\$eruptions,faithful\$waiting,drv=0L,
 degree=0, bandwidth=.5), col = "blue")

with(faithful, plot(eruptions, waiting, cex=2, main = "bandwidth=2", pch="."))
lines(locpoly(faithful\$eruptions,faithful\$waiting,drv=0L,
 degree=0, bandwidth=2), col = "blue")

... Examples



These are smoother than the plots in ELM using <code>base::ksmooth</code>

Applied Statistics I

November 30 2022

Bias and MSE

- Nadaraya-Watson: $\hat{f}_{\lambda}(\mathbf{x}) = \Sigma w_i y_i / \Sigma w_i;$ $w_i = \frac{1}{\lambda} K(\frac{x_i x_o}{\lambda})$
- $\hat{f}_{\lambda}(x)$ is biased

$$E\{\hat{f}_{\lambda}(x)\} \doteq \frac{1}{2}\lambda^{2}f''(x)$$
$$\operatorname{var}\{\hat{f}_{\lambda}(x)\} \doteq \frac{\sigma^{2}}{n\lambda g(x)}\int_{-1}^{1}K^{2}(u)du$$

 $g(\cdot)$ limiting density of x's

SM 10.7.1 (no *n*); ELM 11.1

- could choose λ to minimize MSE = bias² + var, at x
- could choose λ to minimize integrated MSE
- more usual to use cross-validation

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^{n} \{y_i - \hat{f}_{-i}(x_i)\}^2$$

Applied Statistics I November 30 2022

18

Cross-validation

```
library(sm)
hm <- hcv(faithful$eruptions,
   faithful$waiting, display = "lines")
sm.regression(faithful$eruptions,
   faithful$waiting, h = hm,
   xlab = "eruptions",
   ylab = "waiting")</pre>
```



Local Polynomials

- above uses local averaging based on kernel function
- better estimates can be obtained using local regression at point x

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

• attach a weight to each observation (y_i, x_i)

according to the distance from x_i to x_0

$$\hat{\beta} = (X^T W X)^{-1} X^T W y, \quad W = \operatorname{diag}(W), \quad W_i = \frac{1}{\lambda} K(\frac{x_i - x_o}{\lambda})$$

$$\hat{f}_{\lambda}(\mathbf{x}_{\mathsf{o}}) = \hat{\beta}_{\mathsf{o}}$$

• usually evaluate the function at sample points: $\hat{f}_{\lambda}(\mathbf{x}_i), i = 1, \dots, n$

- odd-order polynomials work better than even; usually local linear fits are used
- kernel function is often a Gaussian density, or the tricube kernel ${\cal K}(u)=(1-|u|^3)^3, \quad |u|\leq 1$
- as with N-W (local averaging) estimators, choice of bandwidth controls smoothness
- loess is the most widely used, and is the default in ${\tt ggplot2}$
- fits a local linear regression, but not by least squares
- uses a robust version of least squares that downweights outliers
- the result is that the bandwidth can change with *x*

... local polynomials

- $\hat{\beta} = (X^T W X)^{-1} X^T W y$ $W = diag(w_1, \dots, w_n)$
- $\hat{f}_{\lambda}(\mathbf{x}_{o}) = \hat{\beta}_{o} = \sum_{i=1}^{n} S(\mathbf{x}_{o}; \mathbf{x}_{i}, \lambda) \mathbf{y}_{i}$
- $S(x_0; x_1, \lambda), \dots, S(x_0; x_n, \lambda)$ first row of "hat" matrix
- this makes it relatively easy to analyse the behaviour of local polynomial smoothers

SM 10.7

- and to simplify the expression for the cross-validation criterion $CV(\lambda)$
- fitting at each sample value gives

$$\hat{f}_{\lambda}(\mathbf{x}_i) = \sum_{j=1}^n \mathsf{S}(\mathbf{x}_i; \mathbf{x}_j, \lambda) \mathsf{y}_j$$

Cross-validation

.

.

$$CV(\lambda) = \sum_{i=1}^{n} \{y_i - \hat{f}_{-i}(x_i)\}^2$$

• for local polynomials

$$CV(\lambda) = \sum_{i=1}^{n} \left\{ \frac{y_i - \hat{f}_{\lambda}(x_i)}{1 - S_{ii}(\lambda)} \right\}^2$$

• even simpler

$$GCV(\lambda) = \sum_{i=1}^{n} \left\{ \frac{y_i - \hat{f}_{\lambda}(x_i)}{1 - \operatorname{tr}(S_{\lambda})/n} \right\}^2$$

$$\hat{f}_{\lambda}(\mathbf{x}_i) = \sum_{j=1}^n \mathsf{S}(\mathbf{x}_i; \mathbf{x}_j, \lambda) \mathbf{y}_j$$

Examples



Examples loess



geom_smooth in ggplot uses local polynomial fitting

Applied Statistics I November 30 2022

robustified

Example



10 · Nonlinear Regression Models

Construction of a local linear smoother. Left panel: observations in the shaded part of the panel are weighted using the kernel shown at the foot. with h = 0.8, and the solid straight line is fitted by weighted least squares. The local estimate is the fitted value when $x = x_0$, shown by the vertical line. Two hundred local estimates formed using equi-spaced x_0 were interpolated to give the dotted line, which is the estimate of g(x). Right panel: local linear smoothers with h = 0.2(solid) and h = 5 (dots).

Recall that a kernel function w(u) is a unimodal density function symmetric about u = 0 and with unit variance. One choice of w is the standard normal density. Another is a rescaled form of the *tricube* function

Applied Statistics I November 30 2022
$$w(u) = \begin{cases} (1 - |u|^3)^3, & |u| \le 1, \\ 0, & \text{otherwise,} \end{cases}$$
 (10.37)

Springer Texts in Statistics

Gareth James Daniela Witten Trevor Hastie Robert Tibshirani

An Introduction to Statistical Learning

with Applications in R

Applied Statistics I N

Example



7.6 Local Regression 281

FIGURE 7.9. Local regression illustrated on some simulated data, where the blue curve represents f(x) from which the data were generated, and the light orange curve corresponds to the local regression estimate $\hat{f}(x)$. The orange colored points are local to the target point x_0 , represented by the orange vertical line. The yellow bell-shape superimposed on the plot indicates weights assigned to each point, decreasing to zero with distance from the target point. The fit $\hat{f}(x_0)$ at x_0 is obtained by fitting a weighted linear regression (orange line segment), and using the fitted value at x_0 (orange solid dot) as the estimate $\hat{f}(x_0)$.

Inference after fitting local polynomials

- model: $y_i = f(x_i) + \epsilon_i$, i = 1, ..., n; $E(\epsilon_i) = 0$; $var(\epsilon_i) = \sigma^2$
- $\hat{f}_{\lambda}(\mathbf{x}_{o}) = \hat{\beta}_{o} = \sum_{i=1}^{n} S(\mathbf{x}_{o}; \mathbf{x}_{i}, \lambda) \mathbf{y}_{i}$
- $\mathsf{E}\{\hat{f}_{\lambda}(x_{o})\} =$
- $\operatorname{var}\{\hat{f}_{\lambda}(x_{o})\} =$
- · how many parameters did we fit?
- by analogy with least squares, estimates of 'degrees of freedom' are $\nu_1 = tr(S_{\lambda})$, or $\nu_2 = tr(S_{\lambda}^T S_{\lambda})$

$$\tilde{\sigma}^2 = \frac{1}{n - 2\nu_1 + \nu_2} \sum \{y_i - \hat{f}_\lambda(x_i)\}^2$$

Applied Statistics I November 30 2022

... inference after fitting local polynomials

• $\mathsf{E}\{\hat{f}_{\lambda}(\mathsf{x}_{\mathsf{O}})\} = \sum_{i=1}^{n} \mathsf{S}(\mathsf{x}_{\mathsf{O}};\mathsf{x}_{i},\lambda)f(\mathsf{x}_{i}), \quad \mathsf{var}\{\hat{f}_{\lambda}(\mathsf{x}_{\mathsf{O}})\} = \sigma^{2}\sum_{i=1}^{n} \mathsf{S}^{2}(\mathsf{x}_{\mathsf{O}};\mathsf{x}_{i},\lambda)$

$$\frac{\hat{f}_{\lambda}(x_{0}) - \mathsf{E}\{\hat{f}_{\lambda}(x_{0})\}}{\widehat{\mathsf{var}}\{\hat{f}_{\lambda}(x_{0})\}^{1/2}} \stackrel{.}{\sim} \mathsf{N}(\mathsf{O},\mathsf{1})$$



Applied Statistics I November 30 2022

30

... inference after fitting local polynomials





- model $y_i = f(x_i) + \epsilon_i$ $f(\cdot)$ "flexible"
- above $f(\cdot)$ is estimated at several points using local constants or local linear regression KernSmooth::locpoly
- \cdot another popular approach is to use some very flexible, but parametric form, for f
- for example, $f(x) = \sum_{m=1}^{M} \beta_m \phi_m(x)$
- examples of ϕ_m : 1, x, x^2 , x^3 ; 1, $\sin(x)$, $\cos(x)$, $\sin(2x)$, $\cos(2x)$
- popular choice piecewise polynomials: e.g. knots at $\xi_1,\xi_2\in[0,1]$
- basis functions $\phi(\mathbf{x})$: 1, \mathbf{x} , \mathbf{x}^2 , $= \mathbf{x}^3$, $(\mathbf{x} \xi_1)^3_+$, $(\mathbf{x} \xi_2)^3_+$
- ELM p.219 builds these "by hand"
- splines::bs() builds cubic splines automatically

Multiple R-squared: 0.3426, Adjusted R-squared: 0.33478 F-statistic: 43.777 on 3 and 252 DF, p-value: < 2.22e-16

> model.matrix(examod)

(Interce	pt)	bs(x,	3)1	bs(x,	3)2	bs(x,	3)3
1	1 0.0	000000	e+00 (0.0000000	0000	0.0000000	e+00
2	1 1.1	399500e	e-02 (0.00004381	762	5.61423426	e-08
Applied Statistics I	1 2.0	569481e	e-02 (0.00014401	1768	3.3611431	e-07
	4 0 5	070000	~~ ~				~~~







Smoothing splines

•
$$y_i = f(x_i) + \epsilon_i$$
, $i = 1, ..., n$

• choose $f(\cdot)$ to solve

$$\min_{f}\sum_{i=1}^{n} \{y - f(x_i)\}^2 + \lambda \int_a^b \{f''(t)\}^2 dt, \quad \lambda > \mathbf{0}$$

• solution is a cubic spline, with knots at each observed x_i value

see SM Figure 10.18 for a non-regularized solution

- has an explicit, finite dimensional solution
- $\hat{f} = \{\hat{f}(x_1), \dots, \hat{f}(x_n)\} = (I + \lambda K)^{-1} y$
- K is a symmetric $n \times n$ matrix of rank n 2

Nonparametric regression

• $y_i = f(x_i) + \epsilon_i$

- local polynomial regression stats::loess, KernSmooth::locpoly ELM 11.3; SM 10.7.1
- regression splines splines::bs , splines::ns
- smoothing splines stats: smooth.spline
- penalized splines pspline::smooth.Pspline
- wavelets wavethresh::wd
- and more...

Peng et al. 2006 ELM 11.4 ELM 11.5: ISLR Ch.7

ELM 11.2b p 218ff

ELM 11.2a; SM 10.7.2

- same ideas can be applied to generalized linear models
- replace linear predictor $\eta_i = \mathbf{x}_i^T \beta$ with $f(\mathbf{x}_i)$
- use local poly, reg splines, etc.

SM Ex. 10.32 logistic regression

Example: ELM Exercise 14.5

- > data(aatemp)
- > plot(year ~ temp, aatemp)
- > View(toronto) # not shown
- > plot(year ~ temp, toronto)



... Example: ELM Exercise 14.5



... Example: ELM Exercise 14.5



40

Example: logistic regression

-51	6
	0

10 · Nonlinear Regression Models

City	Rain	r/m	City	Rain	r/m	City	Rain	r/m	City	Rain	r/m
1	1735	2/4	11	2050	7/24	21	1756	2/12	31	1780	8/13
2	1936	3/10	12	1830	0/1	22	1650	0/1	32	1900	3/10
3	2000	1/5	13	1650	15/30	23	2250	8/11	33	1976	1/6
4	1973	3/10	14	2200	4/22	24	1796	41/77	34	2292	23/37
5	1750	2/2	15	2000	0/1	25	1890	24/51			
6	1800	3/5	16	1770	6/11	26	1871	7/16			
7	1750	2/8	17	1920	0/1	27	2063	46/82			
8	2077	7/19	18	1770	33/54	28	2100	9/13			
9	1920	3/6	19	2240	4/9	29	1918	23/43			
10	1800	8/10	20	1620	5/18	30	1834	53/75			

Table 10.19

Toxoplamosis data: rainfall (mm) and the numbers of people testing positive for toxoplasmosis, r, our of m people tested, for 34 cities in El Salvador (Efron, 1986).

Terms	df	Deviance
Constant	33	74.21
Linear	32	74.09
Quadratic	31	74.09
Cubic	30	62.63

 Table 10.20
 Analysis of deviance for polynomial logistic models fitted to the toxoplasmosis data.

Applied Statistics I November 30 2022

 \longrightarrow toxoplasmosis.Rmd

SM Ex.10.29 and 10.32

Example: logistic regression



Figure 10.17 Local fits to the toxoplasmosis data. The left panel shows fitted probabilities $\widehat{\pi}(x)$, 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 = 400and its 0.95 confidence band. Note the increased variability due to the quadratic fit, and its stronger curvature at the boundaries.

