

COPLOTS, NONPARAMETRIC REGRESSION, AND CONDITIONALLY PARAMETRIC FITS

BY WILLIAM S. CLEVELAND

AT & T Bell Laboratories

Conditionally parametric regression surfaces provide parsimonious fits in the common situation where the effects due to some factors are small but significant, and the effects of other factors are large, complicated, and require nonparametric fitting. One approach to nonparametric regression is local regression, specifically the local fitting of linear and quadratic polynomials of the factors. Recent work of Fan and of Hastie and Loader has shown that local regression is superior to kernel estimation and modified kernel estimation, two methods that have had extensive theoretical investigation and that work poorly in practice. Local regression can be modified in a simple way to produce conditionally parametric fits. The coplot is a graphical method that is particularly helpful for carrying out regression studies, in particular, for determining factors that can be taken to be conditionally parametric.

1. Introduction. Graphical methods are critical tools for analyzing and modeling multivariate data. For example, coplots are particularly useful for regression studies, both parametric and nonparametric (Cleveland (1993 to appear)). Coplots are the first topic of the paper.

Many multivariate data sets are better fitted by nonparametric regression surfaces than parametric ones because the latter lack the flexibility to track all but very simple patterns. One method of nonparametric regression is *local fitting*. The method is an old idea first used by time series analysts to smooth their data (Macauley (1931)). It was brought to more general regression studies in the 1970s (Cleveland (1979), Stone (1977)). Making this computer-intensive method a practical one required that a number of computational methods be developed (Cleveland and Grosse (1991)), but this work proved successful, and local regression is now regularly used in regression studies (Cleveland et al. (1991)). There are a multitude of other nonparametric regression methods. One is splines (Reinsch (1967), Silverman (1985), Wahba (1978)). But for two or more factors, splines have resisted solutions to nasty computational problems, and at this writing are still an n^3 operation. Another

AMS 1980 Subject Classifications: 62-09, 62G07.

Key words and phrases: Graphics, loess, smoothing.

is kernel estimates (Nadaraya (1964), Watson (1964)) and their adaptive versions (Gasser and Müller (1979)). But kernel estimates have the disadvantage that they do not work. All it takes is a serious effort to use them in practice to find this out. New and exciting theoretical results of Fan (Fan (1992, 1993)) and Hastie and Loader (Hastie and Loader (1993)) show why this is the case and why local fitting of linear and quadratic polynomials does far better. The second topic in this paper is nonparametric regression, in particular, what works and what does not.

Nonparametric regression surfaces can consume a large number of degrees of freedom in following the patterns of complex data. One way to conserve degrees of freedom is to fit a special class of nonparametric regression surfaces called *conditionally parametric fits*. This class is the third topic of the paper.

Finally, an appendix provides information about a package of public-domain subroutines, written in Fortran and C, that carry out local fitting, including conditionally parametric fitting.

2. Coplots. Figure 1 is a scatterplot matrix that graphs data from an industrial experiment (Brinkman (1981)). A single cylinder engine was run and three variables were measured. The response, which will be denoted by NO_x , is the concentration of NO plus the concentration of NO_2 in the engine exhaust, normalized by the amount of work of the engine. The units are μg of NO_x per joule. One factor is the equivalence ratio, E , at which the engine was run. E is a measure of the richness of the air and fuel mixture; as E increases, the amount of fuel increases. The second factor is C , the compression ratio of the engine. There were 88 runs of the experiment, so the data consist of 88 measurements of three variables. The scatterplot of NO_x against E shows a strong nonlinear dependence with a peak between 0.8 and 1.0. The scatterplot of NO_x against C shows no apparent dependence; however, we should not at this point draw any firm conclusion since it is possible that a dependence is being masked by the strong effect of E . The scatterplot of C and E shows that the values of the two variables are nearly uncorrelated and that C takes on one of five values.

These data have been ill treated in the past. In the original analysis, $\log \text{NO}_x$ was modeled (Brinkman (1981)). This makes the variance of the errors decrease with the mean of the response. A fourth degree polynomial of E and C was fitted to the data by stepwise regression. Such high-order polynomial fitting is at best a dubious practice. Another analysis of these data ignored an interaction between C and E that cannot be removed by transformation (Rodriguez (1985)). Finally, (Gu (1992)) concluded that NO_x does not depend on C , which, as we will see, is manifestly not the case.

Figure 2 is a coplot of NO_x against C given E . The dependence panels are the 3×3 array, and the given panel is at the top. On each dependence panel, NO_x is graphed against C for those observations whose values of E

lie in an interval, and a smooth curve has been added using the nonparametric regression procedure loess, which will be described in the next section. Thus each panel shows how NO_x depends on C for E held fixed to an interval. The intervals are shown on the given panel; as we move from left to right through these intervals, we move from left to right and then from bottom to top through the dependence panels. The intervals have two properties: approximately the same number of observations lie in each interval and approximately the same number of observations lie in two successive intervals. The data analyst specifies the number of intervals, 9 in Figure 2, and the target fraction of points shared by successive intervals, $1/2$ in Figure 2.

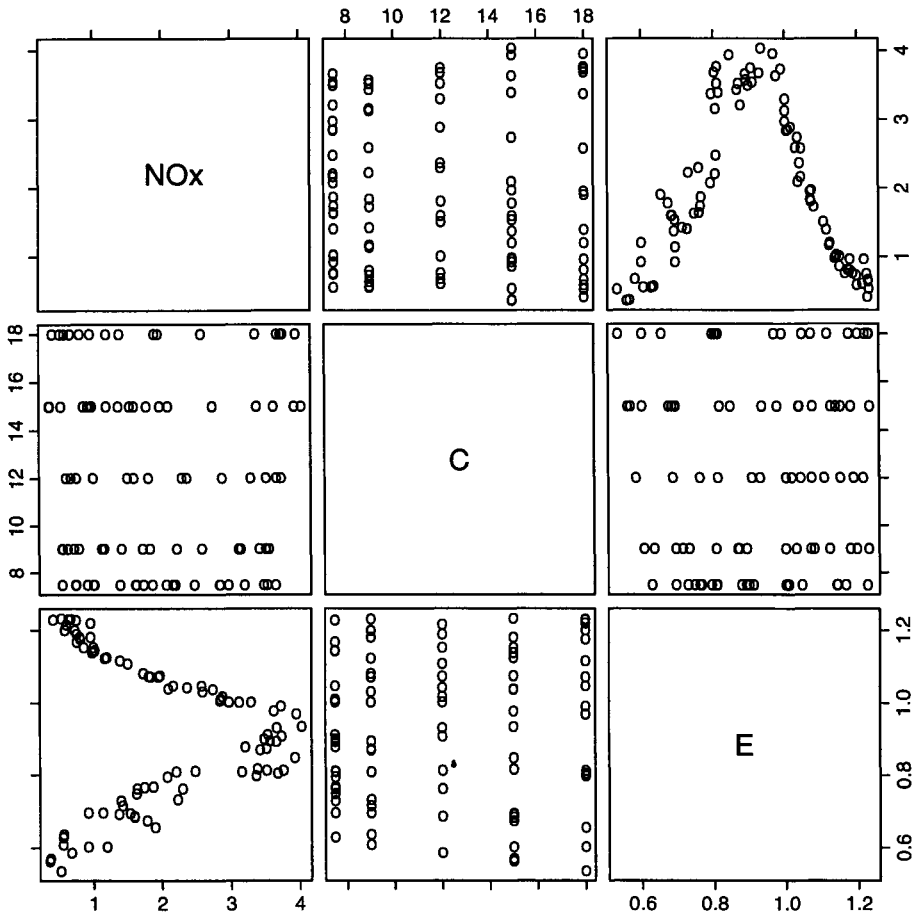


Figure 1. Scatterplot matrix of the engine data

Figure 3 is a coplot of NO_x against E given C . Since C takes on five values, we have simply conditioned on each of these five values.

The coplot in Figure 3 shows that NO_x has a strong nonlinear dependence on E with a peak value near $E = 0.9$ for each conditioning on C . The coplot of Figure 2 shows that NO_x depends on C , but less strongly. That is, over the range of values of E and C in the data set, NO_x undergoes greater change as a

function of E for C held fixed than as a function of C for E held fixed. Given E , the surface appears to be linear as a function of C ; as E increases, the slope first increases and then decreases to zero. The complex nonlinear behavior of the data as a function of E given C makes the fitting of a parametric surface an unsatisfactory approach. Nonparametric regression is the appropriate tool here. But the observation of a simple function of C given E is exactly the property that will be exploited by the conditionally parametric fits of Section 4.

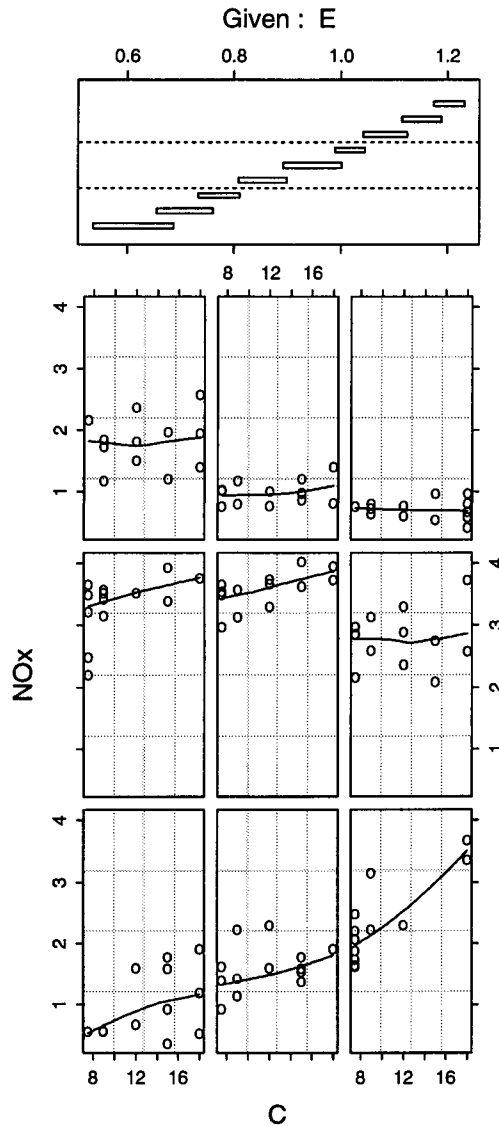


Figure 2. Coplot of NO_x against C given E

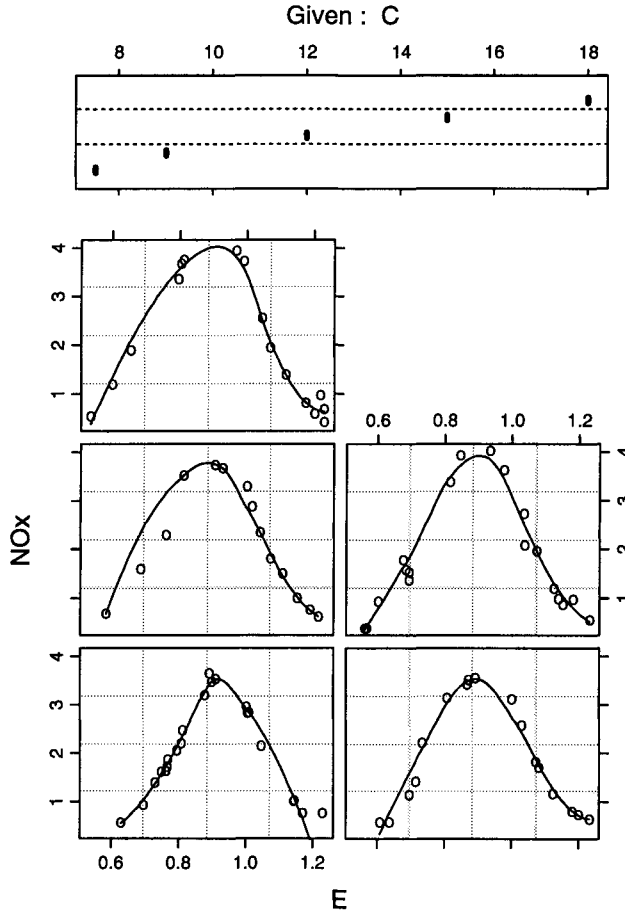


Figure 3. Coplot of NO_x against E given C

3. Nonparametric Regression

3.1. Local Regression Models

Let y_i , for $i = 1$ to n , be a measurement of the response, and let x_i be a vector of measurements of p factors. The model for local regression has the same basic structure as that for parametric regression:

$$y_i = g(x_i) + \varepsilon_i,$$

where g is the regression surface and the ε_i are random errors. In carrying out a local regression, we specify properties of the errors and the regression surface. We will suppose the errors are i.i.d. Gaussian, or relax this assumption and suppose only that they are i.i.d. and symmetric. For each x in the space of the factors, we suppose that in a certain neighborhood of x , the regression surface is well approximated by a linear or quadratic polynomial in the factors. We will let λ be the degree of the polynomial. The overall sizes of the neighborhoods

are specified by a parameter, α , that will be defined shortly. Size, of course, implies a metric in the space of the factors, which we will also define.

3.2. Loess

Loess is one particular method of estimation for local regression models (Cleveland (1979), Cleveland et al. (1988); Cleveland and Devlin (1988)). The following describes the loess computation at one point, x , in the space of the factors.

Let $\Delta_i(x)$ be the Euclidean distance from x to x_i . This provides a metric in the space of the factors, but the x_i do not have to be the raw measurements. Typically, it makes sense to take the x_i to be the raw measurements normalized in some way. Here, we normalize by dividing the factors by their 10% trimmed sample standard deviation. Let $\Delta_{(i)}(x)$ be the values of these distances ordered from smallest to largest.

The smoothness of the loess fit depends on the specification of a neighborhood parameter, $\alpha > 0$. Suppose $\alpha \leq 1$. Let q be equal to αn truncated to an integer. Let

$$T(u) = \begin{cases} (1 - u^3)^3, & \text{for } 0 \leq u < 1 \\ 0, & \text{for } u \geq 1 \end{cases}$$

be the *tricube weight function*. We define a weight for (x_i, y_i) by

$$w_i(x) = T(\Delta_i(x)/\Delta_{(q)}(x)).$$

For $\alpha > 1$, the $w_i(x)$ are defined by the same formula except that $\Delta_{(q)}(x)$ is replaced by $\Delta_{(n)}(x)\alpha^{1/p}$. The $w_i(x)$, which we will call the *neighborhood weights*, decrease or stay constant as x_i increases in distance from x .

If we have specified the surface to be locally well approximated by a linear polynomial — that is, if λ is 1 — then a linear polynomial in the factors is fitted to y_i using weighted least squares with the weights $w_i(x)$; the value of this fitted polynomial at x is the loess fit, $\hat{g}(x)$. If λ is 2, a quadratic is fitted. In fitting quadratics it is also sometimes useful to drop the squares of some or even all of the factors in doing the fitting. An example will be given later.

This loess fitting method applies when the data are Gaussian. Modifications can be made in straightforward ways to produce a robust procedure that accommodates a specification of a symmetric, possibly long-tailed, distribution (Brillinger (1977), Cleveland (1979)). In the Gaussian case, sampling distributions of the estimates have been worked out (Cleveland and Devlin (1988)), but only very rough approximations exist for the symmetric case (Cleveland et al. (1991)).

How do we decide between locally linear and locally quadratic fitting? Locally linear fitting is sufficient when the curvature in the surface is relatively gentle. But if there is substantial curvature, in particular, peaks or valleys,

locally quadratic fitting almost always does better. A locally quadratic surface that fits such data without substantial lack of fit will typically be smoother than a locally linear surface that fits the data. The reason is that the added flexibility of the quadratic polynomial often allows us to substantially increase the value of α .

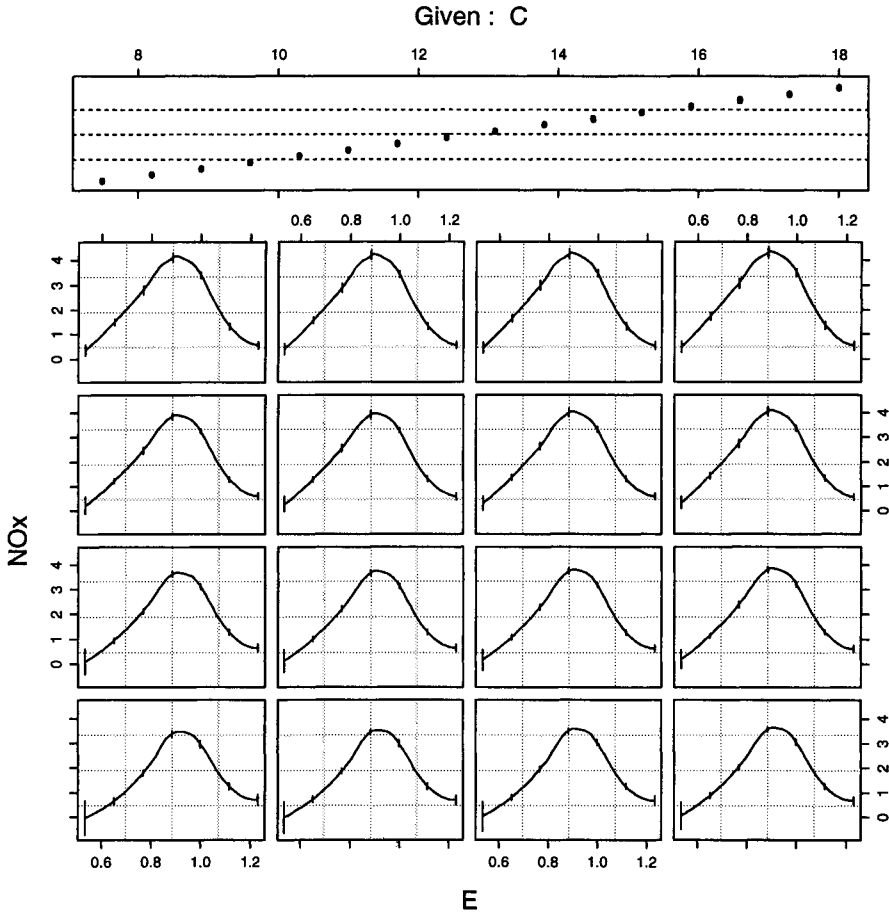


Figure 4. Coplot of local regression fit with pointwise 99% confidence intervals

The engine data described earlier require locally quadratic fitting because there is substantial curvature: a ridge in the surface near $E = 0.9$ that was revealed by the coplots. Diagnostic plotting of residuals shows leptokurtosis, so robust fitting is needed. Diagnostic plotting also shows that with this specification, an α of about $1/4$ is the largest the data will tolerate without lack of fit. The resulting fit has 21.6 equivalent degrees of freedom.

Figure 4 shows a coplot of the surface as a function of E given C . For each of a collection of values of C , shown at the top on the given panel, the surface is graphed as a function of E . In other words, we are graphing slices of the surface using planes perpendicular to the C axis. 99% confidence intervals

are shown at selected positions. Figure 5 shows a coplot of the surface as a function of C given E . There is one distressing property — too much variation as function of C , but we will cure this problem shortly.

3.3. Why Kernel Smoothers Do Not Work

Kernel smoothers amount to locally constant fitting, that is, $\lambda = 0$. In practice, the trade-off between variance and bias almost never makes it worthwhile to use this crude method of approximation. Locally constant fitting cannot even accommodate a linear effect at the boundaries of the design space of the factors. For the NO_x data, a kernel estimate would perform wretchedly. In fact, even locally linear fitting does poorly and locally quadratic fitting is needed.

Fan shows that locally linear estimates have the same asymptotic variance properties as the standard kernel estimates and the same asymptotic bias properties as modified kernel estimates that correct for bias but inflate variances by 50% (Fan (1992, 1993)). In other words, locally linear estimates have the good variance properties of kernel estimates and the good bias properties of corrected-kernel estimates. More recently, an excellent study by Hastie and Loader (Hastie and Loader (1993)) — one that balances mathematics, statistical methods, computing, and the processes of data analysis — provides particularly deep insight into this issue.

4. Conditionally Parametric Fits. A nonparametric surface is conditionally parametric if we can divide the factors up into two disjoint subsets A and B with the following property: given the values of the factors in A , the surface is a member of a parametric class as a function of the the factors in B . We say that the surface is conditionally parametric in A .

It makes sense to specify a regression surface to be conditionally parametric in one or more variables if exploration of the data or a *priori* information suggests that the underlying pattern of the data is globally a very smooth function of the variables. Making such a specification when it is valid can result in a more parsimonious fit.

An exceedingly simple modification of loess fitting yields a conditionally parametric surface. We simply ignore the conditionally parametric factors in computing the Euclidean distances that are used in the definition of the neighborhood weights, $w_i(x)$.

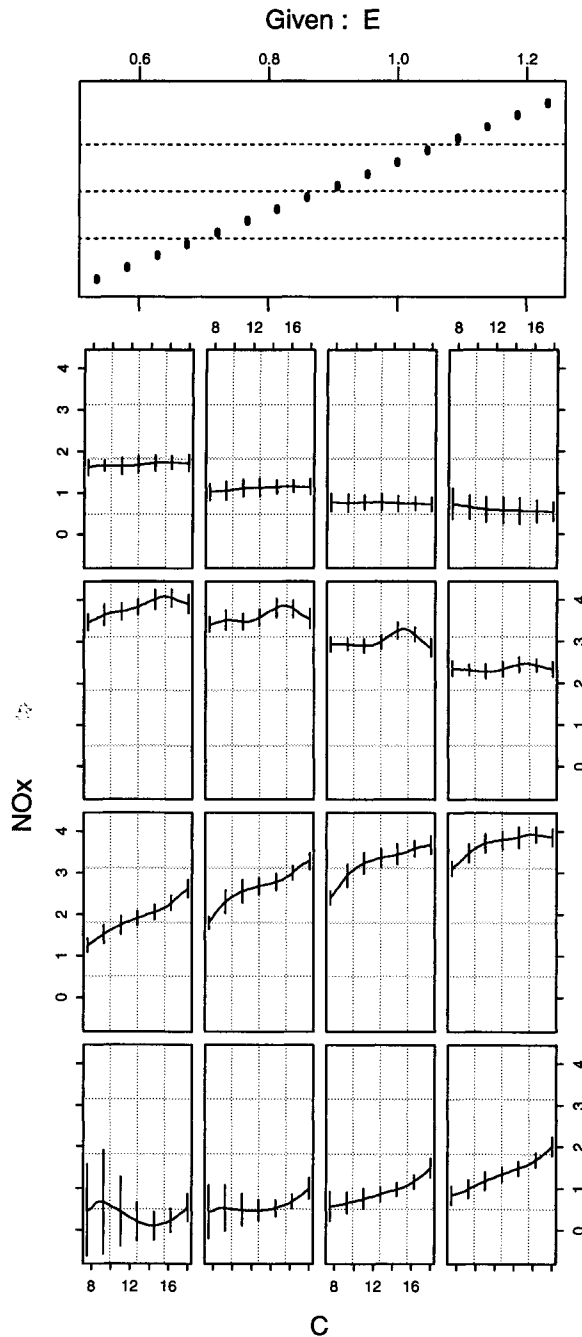


Figure 5. Coplot of local-regression fit with pointwise 99% confidence intervals.

The idea of conditionally parametric surfaces and the approach to fitting them with the simple modification of loess were introduced and put to practical use in (Cleveland et al. (1991)). The idea has been further discussed and expanded in (Hastie and Tibshirani (1993)).

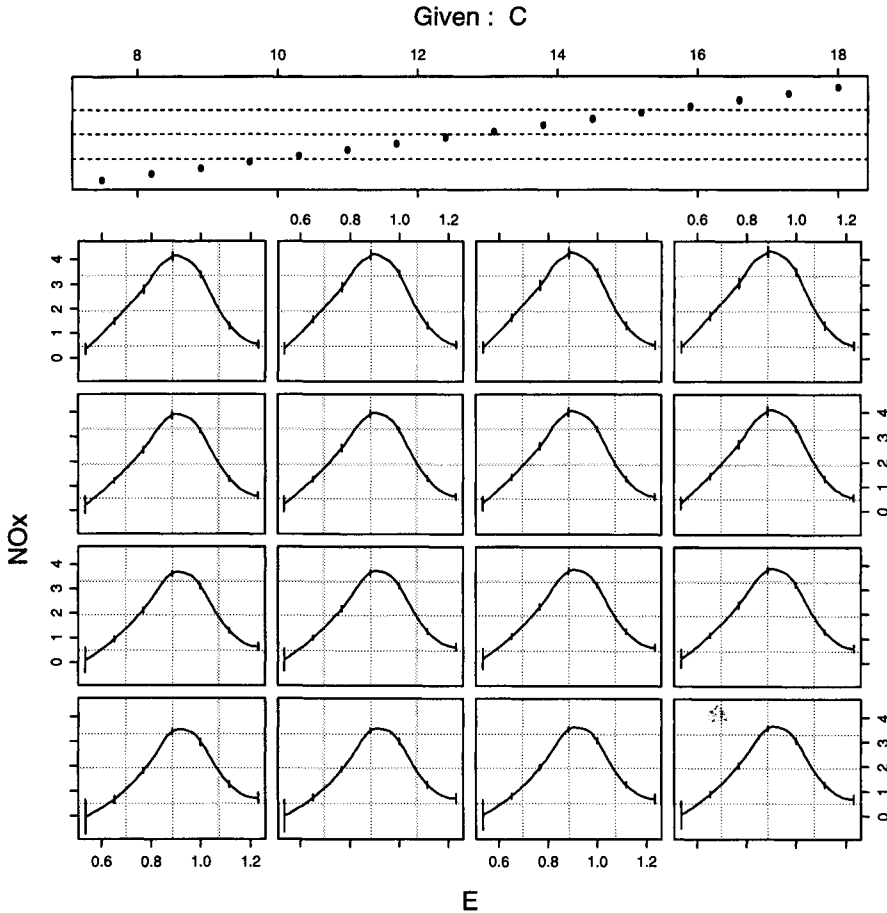


Figure 6. Coplot of conditionally parametric local-regression fit with pointwise 99% confidence intervals.

4.1. Back to the Engine Data

Figure 5, the coplot of the NO_x surface as a function of C given E , does not reflect the linearity in C given E strongly suggested by the coplots. Furthermore, the departures from linearity are not large compared with the sizes of the confidence intervals. We can have conditional linearity in C by specifying the fit to be conditionally parametric in C and dropping C^2 from the local fitting monomials; that is, we fit just C , E , CE , and E^2 . The diagnostic plots of residuals now show that we can increase α to $1/3$ without introducing significant lack of fit. The resulting equivalent degrees of freedom is 13.6, a substantial reduction from the 21.6 of the original fit. Thus by specifying a conditionally parametric fit, we have driven down the degrees of freedom without introducing lack of fit. The coplots in Figures 6 and 7 show the resulting fitted surface.

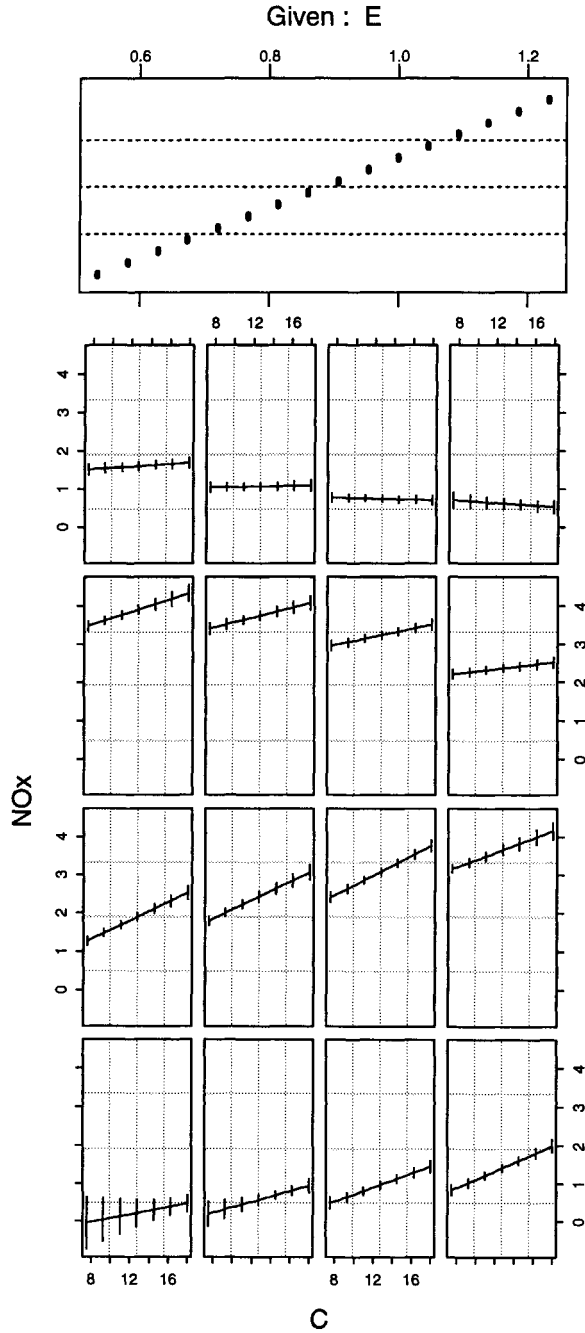


Figure 7. Coplot of conditionally parametric local-regression fit with pointwise 99% confidence intervals.

4.2. Why Conditionally Parametric Fits Work

A number of methods have been introduced to add structure to non-parametric surfaces to cut back on the degrees of freedom of the fit without

sacrificing lack of fit (Breiman and Friedman (1985), Friedman et al. (1983), Friedman and Stuetzle (1981), Hastie and Tibshirani (1990)). In any application, this works if the pattern of the data is reasonably well approximated by such a surface. The above methods ultimately rely on the pattern being well approximated by a sum of univariate functions, either of the original predictors or of linear combinations of them, or of transformations of the response and factors.

A growing list of practical applications has shown that conditionally parametric fits add structure in a way that is quite useful in practice (Cleveland et al. (1991), Hastie and Tibshirani (1993 to appear)). This is to be expected because there is a good argument for why we would expect conditionally parametric surfaces to frequently provide a good fit. If a surface has, instead of no effect, a very limited effect due to these factors compared to other factors, it is reasonable to expect that the effect would be nearly linear or quadratic in these factors given the others. One way to argue for this is to think of Taylor series approximations. Consider the NO_x surface as a function over all imaginable values of C and E . Suppose it is very complex function of both C and E with many peaks and valleys. Taylor series arguments make it clear that C can be varied by a sufficiently small amount in an experiment so that the surface is well approximated by a function that is constant in C given E . If the domain of variation of C is enlarged, the next approximation is a linear function of C given E , and then quadratic. This argument carries with it a corollary; we would expect that, overall, the influence of the variables in which the surface is conditionally parametric would be less than the influence of the other variables. The coplots in Figure 6 and Figure 7 show this to be the case for the engine data; the surface changes less as a function of C , the factor in which the surface is conditionally parametric, than as a function of E .

Acknowledgments. The Symposium on Multivariate Analysis that stimulated this collection of papers was held at Hong Kong Baptist College in March 1992. One extraordinary aspect was its organization. The myriad tasks of the conference were carried out by 78 undergraduates and 4 graduate students in maths and stats. Energy abounded. There was a striking efficiency, enthusiasm, and pride with which they ran the conference. Figure 8 is a quantile plot — the i th order statistic graphed against $(i - 0.5)/78$ — of the time in hours spent by the undergraduates. The mere existence of the data speaks to the fastidiousness of the operation. For a citizen of the United States, even a well-traveled one, this was a revelation. Sometimes the march of history — often thought to move to the loud beat of political leaders' machinations — steps to other things, seemingly small, but the real drivers of revolutions. In this case it is a quest for knowledge, a penchant for perfection, constancy of purpose, and teamwork. Those in the West have felt their effects acutely as drastically shifting national economies. The Symposium provided

the opportunity to see the drivers first hand.

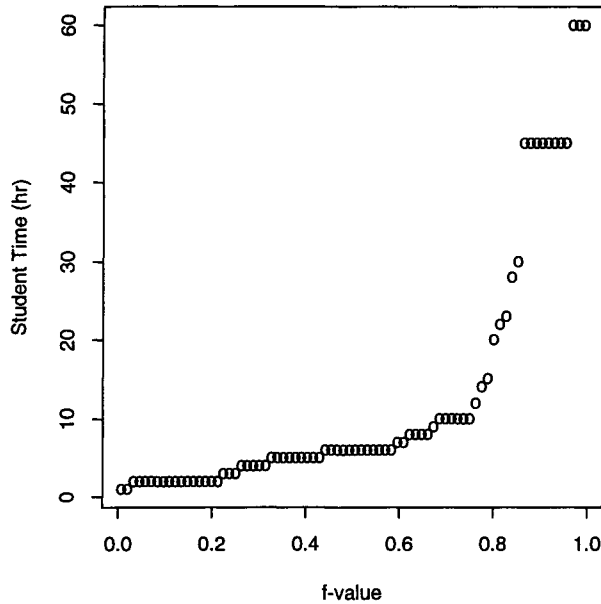


Figure 8. Quantile plot of times spent by undergraduates organizing and running the Symposium.

Appendix: Obtaining the Loess Routines Electronically

The C and Fortran routines, all of which are freely available, may be obtained by sending electronic mail to

NETLIB@RESEARCH.ATT.COM

a mailbox at AT&T Bell Laboratories in Murray Hill, NJ. The message

SEND DLOESS FROM A

should be sent. The routines are double precision.

The file DLOESS is a so-called "shell archive" or "bundle". Moreover, in order to send this 172 kilobyte file to you by email, netlib breaks it into pieces which are themselves shell archives. So you'll need to run SH once on each piece of mail to reconstruct the file DLOESS, then run SH DLOESS to finally reconstruct all the source files.

Subroutines from Linpack, which are called by the Fortran code, are not included. If they are not already on your system, send the message

SEND D1MACH DNRM2 DSVDC DQRDC DDOT DQRSL IDAMAX FROM
LINPACK CORE

to the same address. When installing, don't forget to uncomment the appro-

ropriate DATA statements in dlmach, as described by the comments in those functions.

A PostScript file of a manual on how to use the loess routine for data analysis is also available by email

```
SEND CLOESS.PS FROM A
```

but since it is over half a megabyte, ftp is a better choice

```
FTP RESEARCH.ATT.COM
LOGIN: NETLIB
PASSWORD: {YOUR EMAIL ADDRESS}
BINARY
CD A
GET CLOESS.PS.Z
QUIT
UNCOMPRESS CLOESS.PS
```

Bug reports will receive prompt attention. Send electronic mail to

```
SHYU@RESEARCH.ATT.COM
```

or send paper mail to

```
MING-JEN SHYU
AT&T BELL LABORATORIES
600 MOUNTAIN AVENUE, ROOM 2C-263
MURRAY HILL NJ 07974
USA
```

REFERENCES

- BREIMAN, L. and FRIEDMAN, J.H. (1985). Estimating optimal transformations for correlation and regression. *Journal of the American Statistical Association*, **80**, 580–598.
- BRILLINGER, D. R. (1977). Consistent nonparametric regression. *The Annals of Statistics* **5**, 622–623. Discussion of paper by C. J. Stone.
- BRINKMAN, N. D. (1981). Ethanol fuel — a single-cylinder engine study of efficiency and exhaust emissions. *SAE Transactions*, **90**, 1410–1424.
- CLEVELAND, W. S. (1979). Robust locally-weighted regression and smoothing scatterplots. *Journal of the American Statistical Association*, **74**, 829–836.
- CLEVELAND, W. S. (1993). *Visualizing Data*, Hobart Press, Summit, New Jersey, U.S.A.

- CLEVELAND, W. S., DEVLIN, S. J. and GROSSE, E. (1988). Regression by local fitting: methods, properties, and computational algorithms. *Journal of Econometrics* **37**, 87–114.
- CLEVELAND, W. S. and DEVLIN, S. J. (1988). Locally-weighted regression: an approach to regression analysis by local fitting. *Journal of the American Statistical Association* **83**, 596–610.
- CLEVELAND, W. S. and GROSSE, E. (1991). Computational methods for local regression. *Statistics and Computing*, **1**, 47–62.
- CLEVELAND, W. S., GROSSE, E. and SHYU, W. M. (1991). Local regression models, in *Statistical Models in S*, 309–376, (Eds. J.M. Chambers and T. Hastie). Chapman and Hall, New York.
- FAN, J. (1992). Design-adaptive nonparametric regression. *Journal of the American Statistical Association* **87**, 998–1004.
- FAN, J. (1993). Local linear regression smoothers and their minimax efficiency. *The Annals of Statistics* **21**, 196–216.
- FRIEDMAN, J.H., GROSSE, E.H. and STUETZLE, W. (1983). Multidimensional additive spline approximation. *SIAM Journal of Scientific and Statistical Computing* **4**, 291–301.
- FRIEDMAN, J. H. and Stuetzle, W. (1981). Projection pursuit regression. *Journal of the American Statistical Association* **76**, 817–823.
- GASSER, T. and MÜLLER, H. G. (1979). *Kernel estimation of regression functions*, in *Lecture Notes in Mathematics*, **757**, 23–68, Springer-Verlag, Berlin.
- GU, CHONG. (1992). Diagnostics for nonparametric regression models with additive terms. *Journal of the American Statistical Association* **87**, 1051–1058.
- HASTIE, T. and LOADER, C. (1993). Local regression: automatic kernel carpentry (with discussion). *Statistical Science* **4**, 120–143.
- HASTIE, T. and TIBSHIRANI, R. (1990). *Generalized Additive Models*. Chapman and Hall, London.
- HASTIE, T. and TIBSHIRANI, R., (1993 to appear). Varying-coefficient models (with discussion). *Journal of the Royal Statistical Society*.
- MACAULEY, F.R. 1931). *The Smoothing of Time Series*. New York: National Bureau of Economic Research.
- NADARAYA, E. A. (1964). On estimating regression. *Theory of Probability and its Applications*, **9**, 141–142.
- REINSCH, C. (1967). Smoothing by spline functions. *Numerische Mathematik*

10, 177–183.

- RODRIGUEZ, R. N. (1985). A comparison of the ACE and MORALS algorithms in an application to engine exhaust emissions modeling, in *Computer Science and Statistics: Proceedings of the Sixteenth Symposium on the Interface*, (Ed. L. Billard). North-Holland, New York, 159–167.
- SILVERMAN, B. W. (1985). Some aspects of the spline smoothing approach to nonparametric regression curve fitting (with discussion). *Journal of the Royal Statistical Society B* **47**, 1–52.
- STONE, C. J. (1977). Consistent nonparametric regression. *The Annals of Statistics* **5**, 595–620.
- WAHBA, G. (1978). Improper priors, spline smoothing, and the problem of guarding against model errors in regression. *Journal of the Royal Statistical Society B* **40**, 364–372.
- WATSON, G. S. (1964). Smooth regression analysis, *Sankhya A* **26**, 359–372.

AT&T BELL LABORATORIES
600 MOUNTAIN AVENUE
MURRAY HILL, NJ 07974
USA