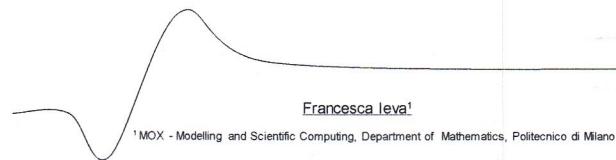


(till now it was testing and CIS, from now it'll be modelling, fitting and prediction)

Nonparametric Regression

We are now
PREDICTION- ORIENTED

Nonparametric Statistics
AA 2020-2021



Outline

- o Why Nonparametric Regression
- o Dataset & Examples
- o Nonparametric Simple Regression
 - Binning and Local Averaging
 - Kernel estimation and Weighted Local Averaging
 - Local Polynomials
 - Splines
 - Regression splines
 - Smoothing splines
- o Nonparametric Multiple regression: GAMs
- o References

Generalized
Additive
Model

2

Outline

- o Why Nonparametric Regression
- o Dataset & Examples
- o Nonparametric Simple Regression
 - Binning and Local Averaging
 - Kernel estimation and Weighted Local Averaging
 - Local Polynomials
 - Splines
 - Regression splines
 - Smoothing splines
- o Nonparametric Multiple regression: GAMs
- o References

3

Why Nonparametric Regression

- o Regression analysis traces the average value of a response variable Y as a function of several predictors $\mathbf{X} = (X_1, \dots, X_p)$
- o The object of regression analysis is to estimate
$$\mu|\mathbf{x} = f(\mathbf{x}) \text{ where } \mu = E[Y]$$
- o Typical (very strong) assumptions:
 - Linear relationship between μ and \mathbf{x} $\Rightarrow f(\mathbf{x})$ is of the form $a\mathbf{x} + b$
 - Omoschedasticity
 - Normal distribution of the errors (then of the responses)
 - Independent observations ($\text{errors } \perp \!\!\! \perp \text{ observations } \perp \!\!\! \perp$)
- o There are many ways in which such assumptions may go wrong
 - In time series errors may not be independent
 - Conditional variance of Y may not be constant (heteroscedasticity)
 - Conditional distribution of Y may not be Normal

Y = output of interest
 \mathbf{x} = input

Remember: when we talk about expected values, we always intend conditional to data ($E[Y|\mathbf{x}]$)

Why Nonparametric Regression

- o Nonparametric regression relaxes the assumption of linearity, substituting the much weaker assumption of smooth regression function f such that

$$Y = f(X) + \varepsilon \quad E[\varepsilon] = 0$$
- o The goal of nonparametric regression is to construct an estimate of f from i.i.d. samples $(y_1, x_1), \dots, (y_n, x_n)$ having the same joint distribution, not assuming any parametric form for f .
- ! Note: we usually assume X as being fixed, i.e. we are conditioning on a given realization of $X=x$, but then errors may result i.i.d. only if we assume independence between errors ε and inputs X .
(and so we make this assumption)
- o Gain -> more accurate estimation => better performance/prediction of the model
- o Cost -> greater computational effort and more difficult-to-understand results (we completely lose, for instance, the interpretability of the betas. Now we won't interpret them)

Outline

- o Why Nonparametric Regression
- o Dataset & Examples
- o Nonparametric Simple Regression
 - Binning and Local Averaging
 - Kernel estimation and Weighted Local Averaging
 - Local Polynomials
 - Splines
 - Regression splines
 - Smoothing splines
- o Nonparametric Multiple regression: GAMs
- o References

6

1. Wage dataset (from JHWT)

- o This application (which is referred to as the Wage dataset throughout the JHWT book) examines a number of factors that relate to wages for a group of males from the Atlantic region of the United States.
- o Goal: to understand the association between an employee's age and education, as well as the calendar year, on his wage.

Wage dataset (from JHWT)

this is the typical case in which a standard (linear) regression would poorly behave.

The other two cases (Year and Education) show that maybe a linear regression would be ok.

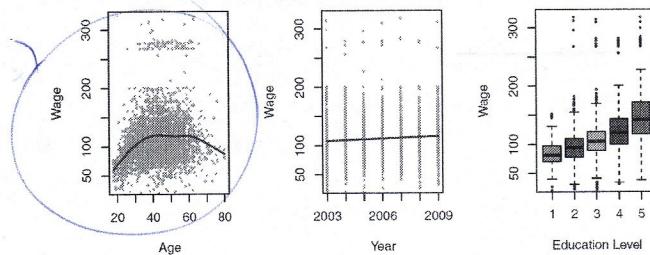


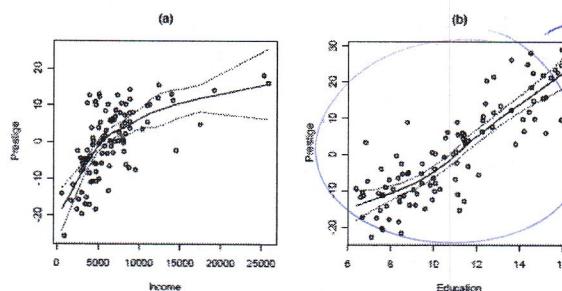
FIGURE 1.1. Wage data, which contains income survey information for males from the central Atlantic region of the United States. Left: wage as a function of age. On average, wage increases with age until about 60 years of age, at which point it begins to decline. Center: wage as a function of year. There is a slow but steady increase of approximately \$10,000 in the average wage between 2003 and 2009. Right: Boxplots displaying wage as a function of education, with 1 indicating the lowest level (no high school diploma) and 5 the highest level (an advanced graduate degree). On average, wage increases with the level of education.

Z. Occupational prestige

- Blishen and McRoberts (1976) reported a linear multiple regression of the rated prestige of 102 Canadian occupations on the income and education levels of these occupations in the 1971 Canadian census.
- Goal: to produce substitute predicated prestige scores for many other occupations for which income and education levels were known, but for which direct prestige ratings were unavailable.
- Figure 3 shows the results of fitting an *additive nonparametric regression* to Blishen's data:

$$y = \alpha + f_1(x_1) + f_2(x_2) + \epsilon$$
- The graphs in Figure 3 show the estimated partial regression functions for income (f_1) and education (f_2).

Occupational prestige



In here we may even say that the relation is linear. However, for the income there is no way to be linear.

Figure 3. Plots of the estimated partial-regression functions for the additive regression of prestige on the income and education levels of 102 occupations.

3. Infant mortality

- Goal: modeling the relationship between infant-mortality rates (infant deaths per 1,000 live births) and GDP per capita (in U.S. dollars) for 193 nations of the world (Figure 1 (a)).
 - Although infant mortality declines with GDP, the relationship between the two variables is highly nonlinear: As GDP increases, infant mortality initially drops steeply, before leveling out at higher levels of GDP.

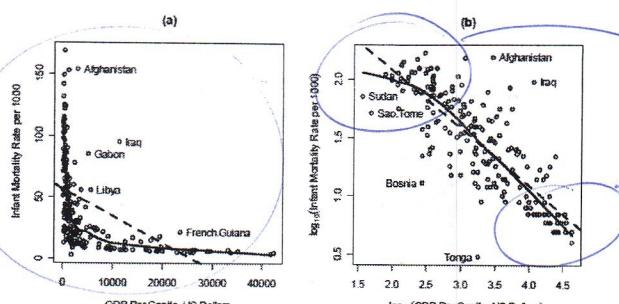
Note: The nonparametric regression line on the graph was produced by a method called *loess* (or *loess*), an implementation of local polynomial regression, and the most commonly available method of nonparametric regression.

- Because both infant mortality and GDP are highly skewed, most of the data congregate in the lower-left corner of the plot, making it difficult to discern the relationship between the two variables.
- The linear least-squares fit to the data does a poor job of describing this relationship.

Note: in Figure 1 (b), both infant mortality and GDP are transformed by taking logs. Now the relationship between the two variables is nearly linear.

Infant mortality

Strongly nonlinear.
The problem is: even if we perform some transit. (e.g. the logarithm) we end up in a situation that is still nonlinear (second picture)



Note: always do a closer look on what's happening on the boundaries: is it a sign of something? can we omit it?

Figure 1. Infant-mortality rate per 1000 and GDP per capita (US dollars) for 193 nations.

Outline

- o Why Nonparametric Regression
- o Dataset & Examples
- o Nonparametric Simple Regression
 - Binning and Local Averaging
 - Kernel estimation and Weighted Local Averaging
 - Local Polynomials
 - Splines
 - Regression splines
 - Smoothing splines
- o Nonparametric Multiple regression: GAMs
- o References

3

Nonparametric Simple Regression

- o Also known as Scatterplot Smoothing *Simple $\Rightarrow \dim(x) = 1$, that's why "scatterplot"* *(the dimension is 1, we can always plot)*
 - o Given a random pair $(Y, X) \in R \times R^p$, the regression function may be expressed as $f(x) = E[Y|X=x]$.
Let consider the case $p=1$ (Simple regression).
 - o The goal of nonparametric regression is to build an estimate of f from $(Y_1, X_1), \dots, (Y_n, X_n)$ having the same joint distribution of (Y, X) , not assuming any parametric form for f .
 - o We will assume X as a fixed/given input, i.e. saying that $Y|X=x$ are i.i.d. implicitly means that we are assuming that input (X) and errors (ε) are independent in the regression model.
- We can build estimates for f essentially in two ways:*
- GLOBAL ESTIMATES
→
f

LOCAL ESTIMATES
- the entire range of X is considered* *only part of the range of X is considered*

Global estimates: POLYNOMIAL REGRESSION

- o Polynomial Regression allows for very flexible curves, describing extremely nonlinear behaviors. *grade of the polynomial (max grade)*
$$y_i = \beta_0 + \sum_{j=1}^d \beta_j x_i^j = \beta_0 + \beta_1 x_i + \dots + \beta_d x_i^d + \varepsilon_i \quad i = 1, \dots, n$$
- o Coefficients can still be estimated through LS.
- o Interpretability of coefficients is no more a point.
- ▲ Being a global approach, individual observation can exert an influence on remote parts of the curve *(this is a disadvantage; all the points influence all the parts in the curve)*

Global estimates: POLYNOMIAL REGRESSION

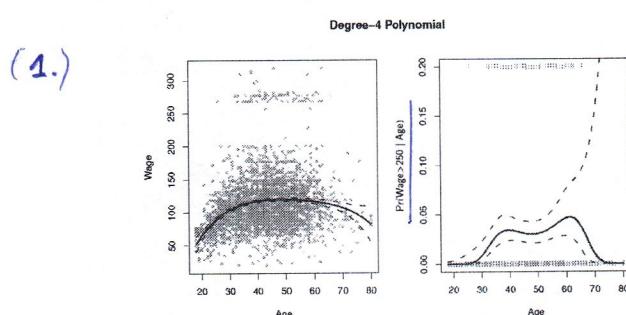


FIGURE 7.1. The Wage data. Left: The solid blue curve is a degree-4 polynomial of wage (in thousands of dollars) as a function of age, fit by least squares. The dotted curves indicate an estimated 95 % confidence interval. Right: We model the binary event $wage > 250$ using logistic regression, again with a degree-4 polynomial. The fitted posterior probability of wage exceeding \$250,000 is shown in blue, along with an estimated 95 % confidence interval.

Global estimates: STEP FUNCTIONS

- Step functions are built selecting K cutpoints, creating K+1 dummy variables indicating the value of the function in each sub interval.

$$y_i = \beta_0 + \beta_1 c_1(x_i) + \dots + \beta_d c_k(x_i) + \varepsilon_i \quad i = 1, \dots, n$$

$$\begin{aligned} c_0(x) &= 1(x < c_1) \\ c_1(x) &= 1(c_1 \leq x < c_2) \\ &\dots \\ c_k(x) &= 1(c_k \leq x) \end{aligned}$$

} indicators

- Coefficients β_j represent the mean in the response for $x \in [c_j, c_{j+1}]$

Global estimates: STEP FUNCTIONS

(1.)

Basically we divide in K portions and in each portion the value is constant (it's the mean of the points belonging to that portion)

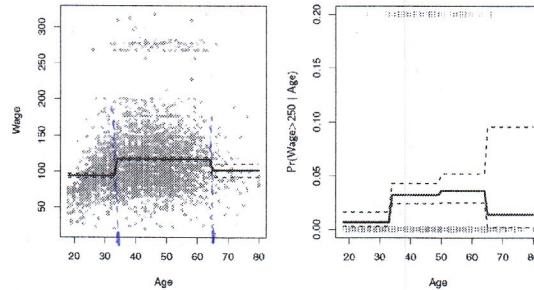


FIGURE 7.2. The Wage data. Left: The solid curve displays the fitted value from a least squares regression of wage (in thousands of dollars) using step functions of age. The dotted curves indicate an estimated 95 % confidence interval. Right: We model the binary event wage > 250 using logistic regression, again using step functions of age. The fitted posterior probability of wage exceeding \$250,000 is shown, along with an estimated 95 % confidence interval.

Towards local estimates: PIECEWISE POLYNOMIALS

- Merging the idea of fitting a polynomial regression and fitting a model on sub intervals of the domain, we get the piecewise polynomial approach.
- In this case, flexibility is guaranteed by fitting low degree polynomials over different regions of X, instead of fitting a high-degree polynomial over the entire range $\Rightarrow \beta_j$ are region specific.
- The points where coefficients change are called knots.

Ex: piecewise cubic polynomial with 1 knot \rightarrow only 2 expressions of Y (before and after the knot)

$$\begin{aligned} y_i &= \beta_{01} + \beta_{11}x_i + \beta_{21}x_i^2 + \beta_{31}x_i^3 + \varepsilon_i & x_i \leq c \\ y_i &= \beta_{02} + \beta_{12}x_i + \beta_{22}x_i^2 + \beta_{32}x_i^3 + \varepsilon_i & x_i > c \end{aligned}$$

- Coefficients can still be estimated through LS.

Note: we don't have to choose the same degree of polynomial in every portion (here we did, but we don't have to)

Towards local estimates: PIECEWISE POLYNOMIALS

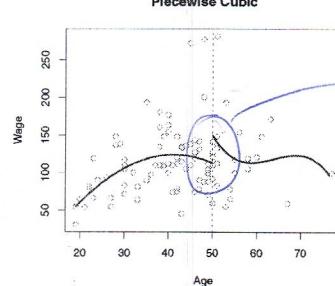
Disadvantages:

- No guarantees about continuity of the fitting.
- High number of parameters to be estimated as long as the number of knots increases.

- Constraints are needed at knots, e.g. continuity on derivatives up to the d-1 degree.
 \rightarrow spoiler: splines

- Note: human eye cannot distinguish discontinuities over second derivatives

here we will ask for the continuity of the derivatives up to the d-1 degree



If we don't impose anything we might end up with something like this, which actually makes no sense! (We want our (global) curve to be at least continuous)

That's why in the literature we often hear about cubic splines

Local estimates: BINNING

- o Suppose that the predictor variable X is continuous.
- o If the sample is large enough, we may divide the range of X in a number of intervals (bins) in order to get the estimate for $\mu|x$.

$$\mu|x \leftarrow \hat{f}(x) = \bar{y}|x$$

we divide the region in bins and we compute the average in each bin

\bar{y}_i = mean of the responses in the i -th bin
(for the inputs belonging in the i -th bin)

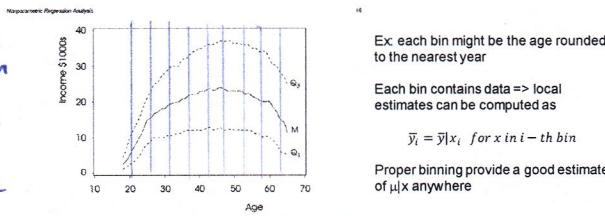


Figure 4. Simple nonparametric regression of income on age, with data from the 1990 U.S. Census one-percent sample.

Local estimates: BINNING

- o Given sufficient data, binning is essentially no-cost.
- o If few observations are present, averages become unstable.
- o A solution may be to allow bins being not uniform in size (as in Figure 5)

--> not always possible to act on bin size.

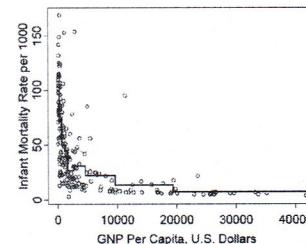


Figure 5. The binning estimator applied to the relationship between infant mortality and GDP per capita.

Local estimates: Bias-Variance trade-off

robustness vs. accuracy
(being close to the data)

Bias-Variance trade-off

$$MSE[\hat{f}(x)] = V[\hat{f}(x)] + Bias^2[\hat{f}(x)]$$

- o Minimizing MSE work at cross purpose:
 - Wide bins produce small variance and large bias in estimates
 - Small bins produce large variance and small bias in estimates

Only with many observations we can care of both jointly (many observations that are well distributed)
- o All nonparametric methods bump up against this problem.
- o Even if the binning estimator is biased, it is consistent as long as the regression function is smooth.

--> "All we need is to shrink the bins toward 0 as the sample size grows up, but shrinking them sufficiently slow so that the number of observations in each bin grows as well".

] we would like to have minuscule bins but with a lot of points in each bin

One step ahead of just making the binning :

Local estimates: LOCAL AVERAGES

- o The basic idea of local averaging is that, as long as the regression line μ is smooth, observations associated to x values near to focal x_0 are more informative about $f(x_0)$.
- ! o Local averaging is different from binning since data are not dissected into a non-overlapping bins, but we move a window continuously over the data, averaging the observations falling in that window.

We then calculate $\hat{f}(x)$ at a number of focal values, usually equally spread.

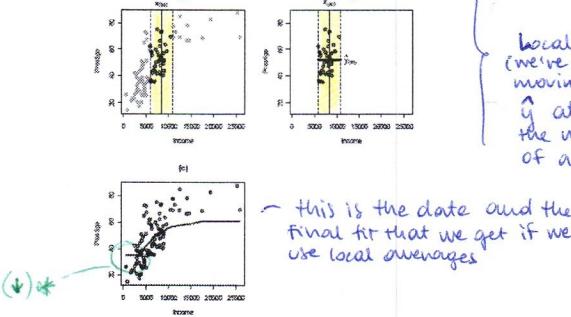
- o Options:

- Fixed width (w) windows, centered at focal values
- Width-adjusted window, containing k points => k Nearest Neighbours (kNN)

] different behavior in bias-variance trade-off

Local estimates: LOCAL AVERAGES

(2.)



local averages (FIXED WINDOW)
 we're continuously moving the window
 at focal point x_0 (the center line of the window) is defined as the mean of all the points contained in the window

→ this is the date and the final fit that we get if we use local averages

Figure 6. Nonparametric regression of prestige on income using local averages.

Local estimates: kNN

we adjust the window every time, forcing it to contain each time k points

- o The kNN estimate of $f(x)$ is

$$\hat{f}(x) = \frac{1}{k} \sum_{i \in N_k(x)} y_i \quad \text{where } N_k(x) \text{ contains the } k \text{ nearest points to } x$$
- o Problems with local averaging :
 - Boundary Bias (artificial flattening of the curve at extremes) (★ (↑))
 - It might produce rough estimates ($\hat{f}(x)$ tends to take small jumps as observations enter/exit the windows)
 - Subject to distortion whenever outliers fall into the window
- o Varying the number of neighbors k we may tune the flexibility of the estimate.
 - The lower the number, the more flexible the estimate (low bias, high variance)
 - The higher the number, the less flexible the estimate (high bias, low variance)
- o The proportion k/n is referred to as span of the estimator.

proportion of data we allow each window to contain

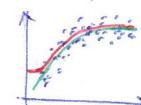
since we want k points, the points that are at the beginning or at the end are influenced only by one-direction:



these points are influenced only by the right-side

these points are influenced only by the left-side

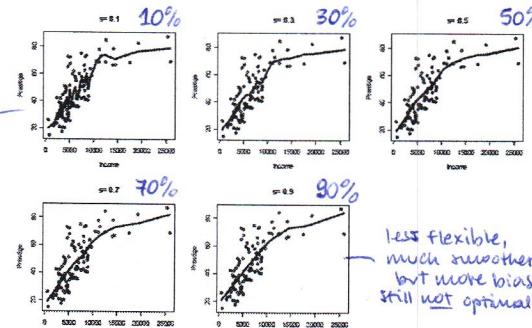
Example :



what we want
 what we get
 (the flattening at the end is not a problem. The flattening at the beginning is)

(2.)

high variability,
 less bias
 (not optimal, the initial part is very "nervous")



less flexible,
 much smoother,
 but more biased,
 still not optimal

made with different span for each fitting
 (the choice of the optimal is something subjective)

Local estimates: WEIGHTED LOCAL AVERAGING

- o If we rewrite

$$\hat{f}(x) = \frac{1}{k} \sum_{i \in N_k(x)} y_i \quad \text{where } N_k(x) \text{ contains the } k \text{ nearest points to } x$$

as $\hat{f}(x) = \sum_{i=1}^n w(x_i) y_i$ where $w(x_i) = \frac{1}{k}$ if x_i is one of the KNN and 0 otherwise

then we may read the estimate as a weighted local averaging.

- o Note that $w(x_i)$ in the case above is a discontinuous function of x
 \Rightarrow the same is true for $\hat{f}(x)$
- o This representation reveals that kNN estimates belong to a class called LINEAR SMOOTHERS.
 Being $y = (y_1, \dots, y_n)$ and $\hat{y} = (\hat{f}(x_1), \dots, \hat{f}(x_n))$, then $\hat{y} = S y$
- o The estimates presented up to now have a really poor dependence on d : as d increases, exponentially more samples are required to achieve a reasonable bound for the MSE
 \Rightarrow Curse of dimensionality!

the weights provide a discontinuous function of x (and so, again, we have to decide the priority: smoothness? accuracy?)

Where S is a linear operator (S is the matrix containing all the weights)

Local estimates: KERNEL ESTIMATION

- Kernel estimation is an extension of local averaging
 \Rightarrow greater weight on observations close to the focal, according to a kernel density
- Let z_i denote the scaled signed distance between the i -th observation and the focal

$$z_i = \frac{x_i - x_0}{h} \quad \xrightarrow{\text{Bandwidth}} \text{scale factor playing the role of window's width in the binning}$$

- A Kernel $K(z)$ is a function attaching greater weights to observations close to x_0 and that falls off symmetrically and smoothly as $|z|$ grows.
- Denoting $w_i = K(\frac{x_i - x_0}{h})$, the kernel estimation at focal point is

$$\hat{f}(x) = \hat{y}|_{x_0} = \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i}$$

It's still a local method;
in the local neighbourhood
we weight data points differently

Local estimates: KERNEL ESTIMATION

- Popular kernel (h fixed) are:
 - Gaussian,
 - Tricube,
 - Rectangular,
 - Epanechnikov

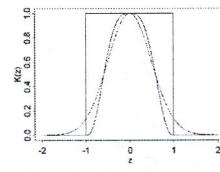


Figure 7. Tricube (light solid line), normal (broken line, rescaled) and rectangular (heavy solid line) kernel functions.

all the points in the bandwidth weight the same

the closer to the focal, the more weight a point has

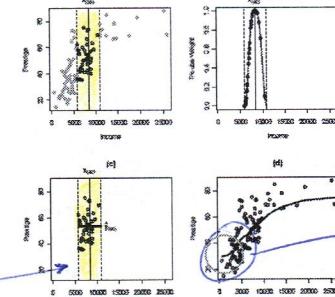
- Also the kernel bandwidth may be made adaptive to contain elements:

$$\frac{k}{n} = \text{SPAN of the kernel}$$

- Varying the bandwidth of the kernel controls the bias-variance trade-off: larger bandwidth \Rightarrow smoother results

Local estimates: KERNEL ESTIMATION

(2.)



in this window we apply a kernel

(in this case it's a bell-shape kernel) and we get an estimate for the focal point that weights the surrounding K -points according to the kernel

Figure 8. The kernel estimator applied to the Canadian occupational prestige data.

the boundary bias is still present (especially for those parts of the curve where we have high derivatives)

Local estimates: LOCAL POLYNOMIALS

- Local polynomials extend kernel estimation to a polynomial fit at the focal point using local kernel weights.

- The resulting Weighted Least Squares (WLS) regression equation is

$$y_i = \beta_0 + \beta_1(x_i - x_0) + \beta_2(x_i - x_0)^2 + \dots + \beta_d(x_i - x_0)^d + \varepsilon_i$$

where the estimates can be obtained minimizing the Weighted Residual Sum of Squares

$$WRSS = \sum_{i=1}^n w_i e_i^2 = \sum_{i=1}^n K\left(\frac{x_i - x_0}{h}\right) (y_i - \beta_0 - \beta_1(x_i - x_0) - \dots)^2$$

- Truncating at $d=1$ we get a local linear regression.
- Truncations at $d>1$ provide more flexibility (i.e., greater variability)
 \Rightarrow Useful to mitigate the boundary bias issue.

Now we have:

- different weights on the observations in the window
- we assume a model for the fitting of the points in each bin
- we average the estimates through a kernel

The result of all of this is a weighted least square regression equation where the estimation of the parameters is weighted through a kernel.

Local estimates: LOCAL POLYNOMIALS

- o The bandwidth h may be fixed or may vary as a function of the focal point, in order to include a given amount of points in the neighbor.
- o When the bandwidth define a window of neighbor (as in the tricube case), the degree of smoothness may be targeted specifying the proportion of observations to be included into the window

$$m = \# \text{ of obs included} = [sn]$$

- o How can we choose the "right" span?
 - Trial and error => start with 0.5 and increase/decrease in order to choose the smallest s providing a smooth fit
 - Cross validation => omit the i -th observation from local regression at the i -th focal value, then evaluate the $CV(s)$ function at different s , selecting the s that act as a minimizer.

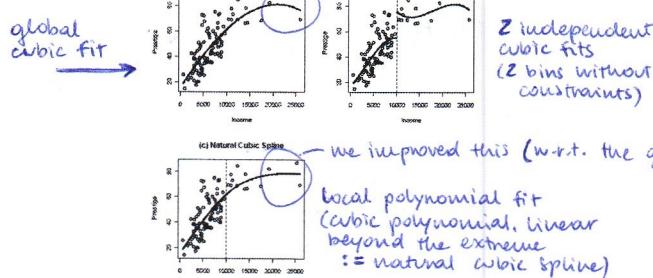
$$CV(s) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i(s)|^2 \quad \text{being } \hat{y}_i(s) = \hat{y}_i |x_i| \text{ for span } s$$

- o Also the local polynomial regression can be seen as a case of linear smoothing.

↑
SPAN of the local regression smoother

Local estimates: LOCAL POLYNOMIALS

(2.)



2 independent cubic fits
(2 bins without constraints)

local polynomial fit
(cubic polynomial, linear beyond the extreme
:= natural cubic spline)

Figure 15. Polynomial fits to the Canadian occupational prestige data: (a) a global cubic fit; (b) independent cubic fits in two bins, divided at Income = 10,000; (c) a natural cubic spline, with one knot at Income = 10,000.

Global

Local estimates: SPLINES

- o Regression and smoothing splines are motivated by a different perspective than kernel and local polynomials.
 - Step back to global estimate, starting from a set of selected basis functions
 - Start from local averaging and move towards higher order local models
- o Splines are piecewise polynomial functions which are constrained to join smoothly at knots.
- o The traditional application for splines is interpolation, but they may be used also for nonparametric regression purposes.
- o Among others, cubic splines are the most popular.
 - (we can impose the smoothness at knots up to the 2nd derivative)

we select suitable basis functions that allow us to fit in each bin a polynomial

This is an upgrade, with the previous method of local polyn. we had the problem of continuity (smoothness in general) at the knots.

Modeling through basis functions guarantees the smoothness at the knots.

We saw:

- general global polynomial fits: able to capture highly varying terms; however they have a high variability, especially at the borders
- local polynomial: we fit data into bins and then we fit polynomial regressors. This approach has the problem of being (possibly) discontinuous.

Solution to the problems:

REGRESSION SPLINES

- o Regression splines guarantee a local fit of polynomials with smooth connections at knots.
 - Cubic regression splines => 3rd order polynomial in each bin with 1st and 2nd derivatives continuous at knots
 - Natural Cubic Splines (NCS) => add knots at boundary and impose linear fit beyond
- o Regression spline "bricks":
 - Outputs: Inputs: y_1, \dots, y_n
 - Inputs: x_1, \dots, x_n
 - # of knots: m
 - Degree of the spline: k
- o The way for ensuring the constraint of smoothness at knots for the function and the $k-1$ derivatives is to build a regression through suitable basis functions

so we'll use k-order splines and they'll be fitted to m selected knots

(which averages between the two) it guarantees the local fitting and the smoothness at the knots.

Local estimates: REGRESSION SPLINES

- o General global polynomial fits are able to capture widely varying forms, but are highly nonlocal and subject to considerable sampling variation.
- o As an alternative, we can fit data into bins, and fit polynomial regression in each bin. --> PB: discontinuity at knots
- o Regression splines guarantee a local fit polynomials with smooth connections at knots.

We're in the univariate case: n observations of the output and the inputs

[Given inputs x_1, \dots, x_n and outputs y_1, \dots, y_n , a k order spline with knots t_1, \dots, t_m may be expressed in the basis form

$$f(x) = \beta_0 + \sum_j \beta_j g_j(x) + \varepsilon$$

$g_j(\cdot)$: basis functions

being the basis composed by a k -degree polynomial and by a truncated power basis function per each knot $\Rightarrow k+m+1$ parameters

$$\Rightarrow x, x^2, \dots, x^k, h(x, t_1), \dots, h(x, t_m)$$

where $h(x, t_j) = (x - t_j)_+^j = (x - t_j)^j$ if $x > t_j$, 0 otherwise

Local estimates: REGRESSION SPLINES

We define

$G \in \mathbb{R}^{n \times (m+k+1)}$,

$$G = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 & \dots & x_1^k & h(x_1, t_1) & \dots & h(x_1, t_m) \\ \vdots & \vdots & \vdots & \vdots & & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & x_n^3 & \dots & x_n^k & h(x_n, t_1) & \dots & h(x_n, t_m) \end{bmatrix}$$

$$\begin{aligned} \hat{\beta} &= \arg \min_{\beta} \|y - G\beta\|_2^2 = \arg \min \sum (y_i - (G\beta)_i)^2 \\ \Rightarrow \hat{\beta}^T &= (G^T G)^{-1} G^T y \\ \Rightarrow \hat{y} &= \hat{f}(x) = G \hat{\beta} \\ &= G(G^T G)^{-1} G^T y \\ &\stackrel{S_{\hat{f}}}{=} \end{aligned}$$

adds the
regression splines
are linear smoother

3rd degree polynomial

Ex: cubic splines

$$f(x) = \beta_0 + (\beta_1 x + \beta_2 x^2 + \beta_3 x^3) + \beta_4 (x - t_1)^3 + \dots + \beta_{3+m} (x - t_{3+m})^3 + \varepsilon$$

Local estimates: REGRESSION SPLINES

- o How can we choose the location and number of knots?

• Location:

regression splines are more flexible in regions containing many knots, since coefficients may change rapidly
 \Rightarrow put more knots where you know the function varies most rapidly, less where is more stable.

• Number:

1) Trial and error 2) CV

(“ignorant choice”)

otherwise, if we know nothing about the problem we can choose to put the points over a grid (equally distanced points)

we capture the same variability that can be captured by high-degree polynomial regression but with fewer parameters (since we do not modify the degree of polynomials, we simply add knots)

Local estimates: REGRESSION SPLINES

- ⚠ Problem: regression splines tend to have high variance at the boundary of the domains

\rightarrow Differently from kernel smoothing, which has poor bias at the boundary

- o Solution: force piecewise polynomial function to have a lower degree to the left of the leftmost knot and to the right of the rightmost knot

\Rightarrow Natural Splines (!!! Only defined for odd order k !!)

- Polynomial of degree k in each interval $[t_1, t_2], \dots, [t_{m-1}, t_m]$
- Polynomial of degree $\frac{k-1}{2}$ on $(-\infty, t_1]$ and $[t_m, \infty)$
- Continuous and with continuous derivatives of order $1, \dots, k-1$ at knots t_1, \dots, t_m

Cubic splines:

in the boundaries the fitting is linear while the fitting in the internal intervals is of degree 3

Local estimates: SMOOTHING SPLINES

- o The idea underlying regression splines relies on specification of a set of knots, producing sequence of basis functions and then using LS for estimating coefficients.
 - o Now we introduce a different approach, also producing splines-based estimates.
 - o Smoothing splines perform a penalized regression over the natural spline basis (orthonormal), placing knots at all the inputs x_1, \dots, x_n
 - => No problem of node selection
 - => Control for overfitting by shrinking the coefficients of the estimated function.
(more shrinkage is assigned to eigenvectors of the basis that correspond to larger eigenvalues)
 - o Smoothing splines can be motivated also from a functional minimization perspective.
-

Local estimates: SMOOTHING SPLINES

- o In fitting a smooth curve to a set of data, we look for a function $g(x)$ minimizing

$$RSS = \sum_{i=1}^n (y_i - g(x_i))^2$$
 - o Note that if we don't constrain $g(x)$, we can always get RSS=0 simply interpolating the data (overfit).
 - => We want a smooth function $g(x)$ making RSS to be small.
 - o A good fit to the data is not only aiming at curve fitting, but also displaying a not "too rapid fluctuation".
 - => Reformulate the problem as a roughness penalty approach, quantifying the notion of "rapidly fluctuating curve" and posing the estimation problem in a way that makes it explicit the necessary trade-off between conflicting aims.
-

Local estimates: SMOOTHING SPLINES

How can we quantify roughness?

- o Given a curve g defined on $[a, b]$, there are many different ways of measuring how rough g is. One may be

$$\int_a^b (g''(t))^2 dt$$

see also [3, p 11] for a more general setting

("rough" function \simeq funzioni "nervose")

- o Why should be this reasonable and/or effective?
 - For not being affected by addition of constant or linear functions
 - $|g''|$ is the maximum number of inflection points of the curve g , but $\int_a^b (g''(t))^2 dt$ is a global measure of roughness
 - Computational advantages
-

Local estimates: SMOOTHING SPLINES

Loss+Penalty formulation:

- o Given a twice-differentiable function g defined in $[a, b]$ and a parameter $\gamma > 0$, we may define the cost function

$$S(g) = \sum_{i=1}^n [y_i - g(x_i)]^2 + \gamma \int_a^b (g''(t))^2 dt$$

RSS Roughness Penalty
- o The penalized LS estimator \hat{g} is the minimizer of $S(g)$ over the class of twice-differentiable functions.
- o Natural cubic splines can be demonstrated to result the interpolant with minimal roughness (see [1], p.18).
This comes as a mathematical consequence of choosing $\int_a^b (g''(t))^2$ as a roughness penalty.

→ tell us about how good the fitness is

→ $\gamma \cdot L^2$ penalization over the domain of the second derivative

→ γ has the role of tuning param.

Local estimates: SMOOTHING SPLINES

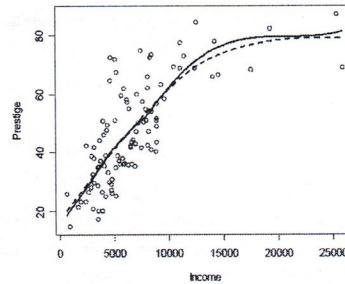
- o The penalty term γ controls the bias-variance trade-off in the smoothing spline
=> the higher γ , the smoother g <=
- o The penalty term ensures that the cost $S(g)$ of a particular curve is determined by the balance of the loss in fitting and penalty on roughness
 - Large γ => RSS is the main component, the minimizer will show small curvature
 $\gamma = \infty$ means high penalty on roughness, making g be the linear regression fit
(high bias, low variance/fluctuating curve)
 - Small γ => Penalty matters, the curve will follow the data closely
 $\gamma = 0$ means no penalty on roughness
(low/no bias, high variance/very fluctuating curve interpolating data)

Local estimates: SMOOTHING SPLINES

- o How can we choose $\gamma \in [0, \infty)$?
Cross validation approach based on LOOCV (Leave One Out Cross Validation) ERROR --> reasonable computational costs
- o Knowing that \hat{g} is a natural cubic spline (NCS) is an enormous advantage, since we need to minimize $S(g)$ over a finite-dimensional class of functions, i.e. the (NCS) with knots at x_i , instead of considering the infinite dimensional set of smooth functions on $[a, b]$.
Optimizing algorithms in [1], pp 18-24.
- o None of the optimizing algorithm depend on the interval $[a, b]$.

Local estimates: SMOOTHING SPLINES

(2.)



Notice that: higher the density of the points in a particular zone and smaller the difference between local linear regression and smoothing (or in general regr.) splines

Remember: splines are the most used.
(all along the nonlinear regression)
local average or something else is used
only if there are not enough data.

Figure 16. Nonparametric regression of occupational prestige on income, using local linear regression (solid line) and a smoothing spline (broken line), both with 4.3 equivalent parameters.

Outline

- o Why Nonparametric Regression
- o Dataset & Examples
- o Nonparametric Simple Regression
 - Binning and Local Averaging
 - Kernel estimation and Weighted Local Averaging
 - Local Polynomials
 - Splines
 - Regression splines
 - Smoothing splines
- o Nonparametric Multiple regression: GAMs
- o References

The multiple regression case: GAMs

- o So far we looked mostly at the **univariate case** $p=1$.
- o Several methods have a **multivariate counterpart** ($p>1$):
 - Kernel smoothing --> multiple dimensional kernel
 - Splines --> thin-plate splines or tensor-product splines [6]
 - Mercer Kernels --> defined directly in multiple dimensions
- o In high dimensions (large p) things are problematic, and suffer from poor variance
- ⇒ Curse of dimensionality, i.e., the estimation gets exponentially harder as the number of dimensions increase.
- o To mitigate this problem, a solution is to adopt an **additive approach**
- ⇒ **GAMs – Generalized Additive Models**
 - Provide a general framework for extending a standard linear model by allowing nonlinear functions of each predictor while maintaining additivity

The multiple regression case: GAMs

- o Instead of considering a full p -dimensional function

$$f(\mathbf{x}) = f(x_1, \dots, x_p),$$
 we restrict our attention to functions of the form

$$f(\mathbf{x}) = f(x_1) + \dots + f(x_p)$$
- o This choice need not to be regarded as an assumption we make about the true function, just like we do not assume that the true model is linear in a traditional regression.
- o Additive models may represent a **good approximation to the truth**, that is able to **scale well with the number of dimensions p** .
- o Stone (1985) shows that while it may be difficult to estimate an arbitrary regression function in high dimension, we can still estimate its best additive approximation well. In other words, each time we cannot hope to recover the true underlying function arbitrarily, we can recover its major structure along the coordinate axes.

The multiple regression case: GAMs

- o **GAMs rely on a sort of building blocks model**, where all the models adopted for each predictor may be packed into one big regression matrix.
- o Estimation with additive models is actually simple: we can just choose our favourite univariate smoother (i.e., nonparametric procedure) and cycle through estimating each function individually (**backfitting**) [3].
 - + GAMs allow for fitting nonlinear models to each predictor
=> we can model relationships that linear models would miss
 - + Due to additivity, we can still analyze and interpret the effect of each covariate on the response, holding all the other covariates fixed
 - GAMs are restricted to be additive

One at the time:
once we're dealing with one variable the others remain fixed.
We update the variable we're dealing with, we fix it and we pass to another variable.

The multiple regression case: GAMs

(1.)

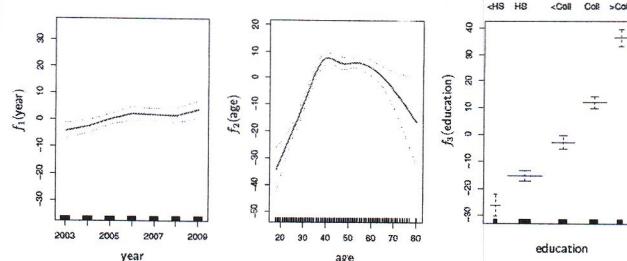


FIGURE 7.11. For the Wage data, plots of the relationship between each feature and the response, wage, in the fitted model (7.16). Each plot displays the fitted function and pointwise standard errors. The first two functions are natural splines in year and age, with four and five degrees of freedom, respectively. The third function is a step function, fit to the qualitative variable education.

The models are different from variable to variable (even between the two continuous): we decide the best fitting separately.

Note: this is different from 3 simple (single) regressions (even if the other covariates are fixed we're considering them)

References

Books & journal papers

- [1] Green, P.J., Silvermann, B.W. (1994) Nonparametric Regression and Generalized Linear Models. A roughness penalty approach. Chapman and Hall
- [2] James, G., Witten, D., Hastie, T., Tibshirani, R. (2013) An Introduction to Statistical Learning with applications in R. Springer
- [3] Tibshirani, R., Wasserman, S. (2015) Statistical Machine Learning course.
- [4] Fox, Nonparametric Simple Regression, Sage (2000)
- [5] Fox, Multiple and Generalized Nonparametric Regression, Sage (2000)
- [6] Wood, S.N. (2003) Thin plate regression splines. Journal of the Royal Statistical Society – series B, 65(1): 95-114

Links

- o <https://statisticallearning.org/nonparametric-regression.html>
- o <https://www.semanticscholar.org/paper/Nonparametric-Regression-Statistical-Machine-%2C-2019-Tibshirani-Wasserman/afb4376dd9983371c648852e8c3ede0cd7666305>
- o <http://www.stat.cmu.edu/~ryantibs/statml/>
- o <https://www.publichealth.columbia.edu/research/population-health-methods/thin-plate-spline-regression>

Packages

splines - <https://www.rdocumentation.org/packages/splines/versions/3.6.2>
gam and mgcv - <https://www.rdocumentation.org/packages/mgcv/versions/1.8-33/topics/gam>

+ Examples:

<https://pjbardelein.github.io/GeogDataAnalysis/sec14.html>