

# Régression polynomiale

2001-2006

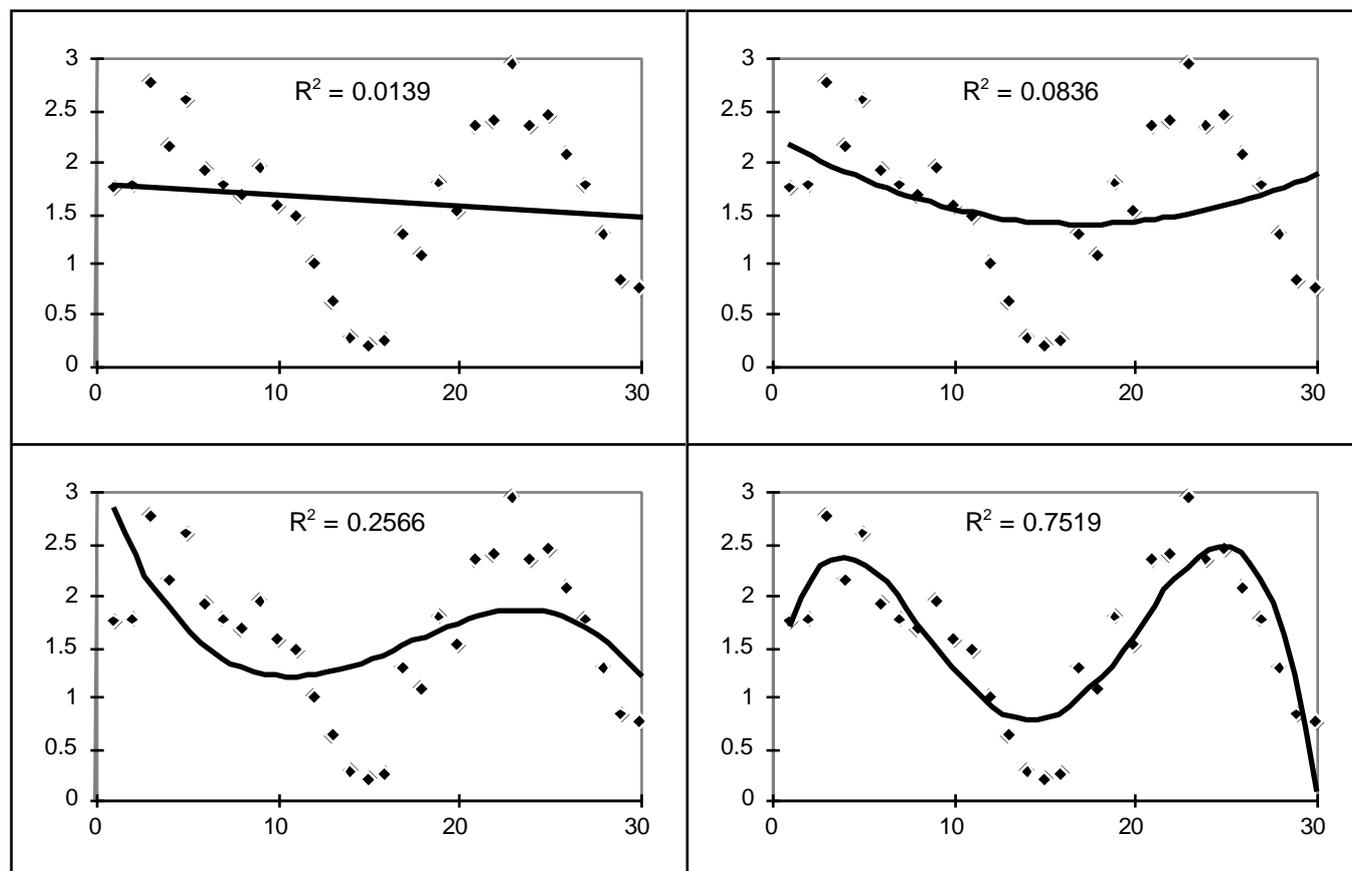
Daniel Borcard, Dép. de sciences biologiques, Université de Montréal

Référence: Legendre et Legendre (1998) p. 526

Une variante de la régression multiple peut quelquefois être appliquée pour ajuster une variable explicative  $x$  (ou plusieurs variables explicatives  $x_j$ ) à une variable dépendante  $y$  de manière non-linéaire. Cette méthode consiste à ajouter à la variable explicative  $x$  de nouvelles variables construites en mettant  $x$  au carré, au cube, etc. L'équation devient donc, pour un polynôme du  $k$ -ième ordre:

$$\hat{y} = a_1x + a_2x^2 + a_3x^3 + \dots + a_kx^k + b$$

Une façon simple de se représenter l'effet de l'ajout d'un ordre est la suivante: **chaque nouvel ordre permet d'ajouter un pli à la courbe**. Une équation du premier ordre est celle d'une droite, du deuxième ordre une parabole; le troisième ordre donne un S couché, etc.:

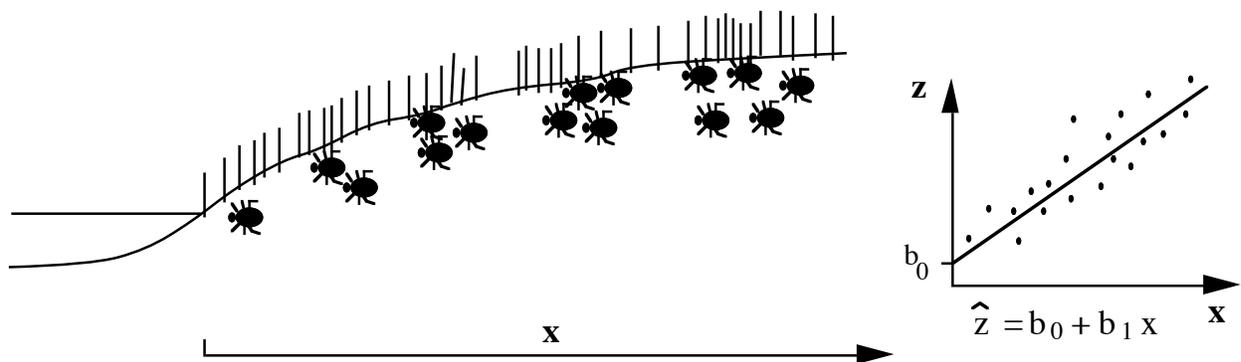


### 6.3.2 Trend surface analysis

This technique is a **particular case of multiple regression**, where the explanatory variables are geographical (x-y) coordinates, sometimes completed by higher order polynomials. When applying this method, one generally supposes that the spatial structure of the observed variable is a result of one or two generating processes that spread over the whole studied area, and that the resulting broad-scale structure of the dependent variable can be modelled by means of a polynomial of the spatial coordinates of the samples. A simple example follows:

Imagine a soil arthropod, the density of which (let us call it  $z$ ) increases from 0 (near a stream) to 100 individuals per square meter (in a nearby meadow). If this density variation is linear, a simple linear regression, with the distance to the stream ( $x$ ) acting as explanatory variable, is enough to model the arthropod density in the whole meadow (Figure 44):

$$\hat{z} = b_0 + b_1 x$$



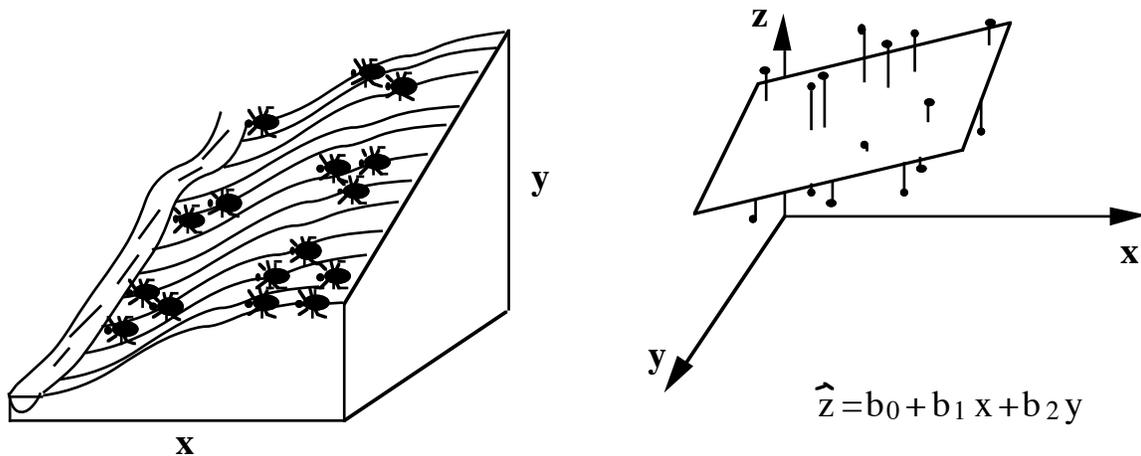
**Figure 44** - Density of an arthropod species along a gradient and linear model.

Now, if the stream (with its neighbouring meadows) extends from higher mountains to sea level, perhaps the arthropod density varies also with the altitude ( $y$ ). A second explanatory variable is necessary, i.e. the altitude, or possibly the distance to the source

along the stream. If the density variation with respect to the altitude is also linear, one gets a first order multiple regression equation of the form:

$$\hat{z} = b_0 + b_1x + b_2y$$

The result is thus a *regression plane* fitted through the  $z$  data (densities) by means of the  $x$ - $y$  coordinates of the arthropod sampling points (Figure 45).

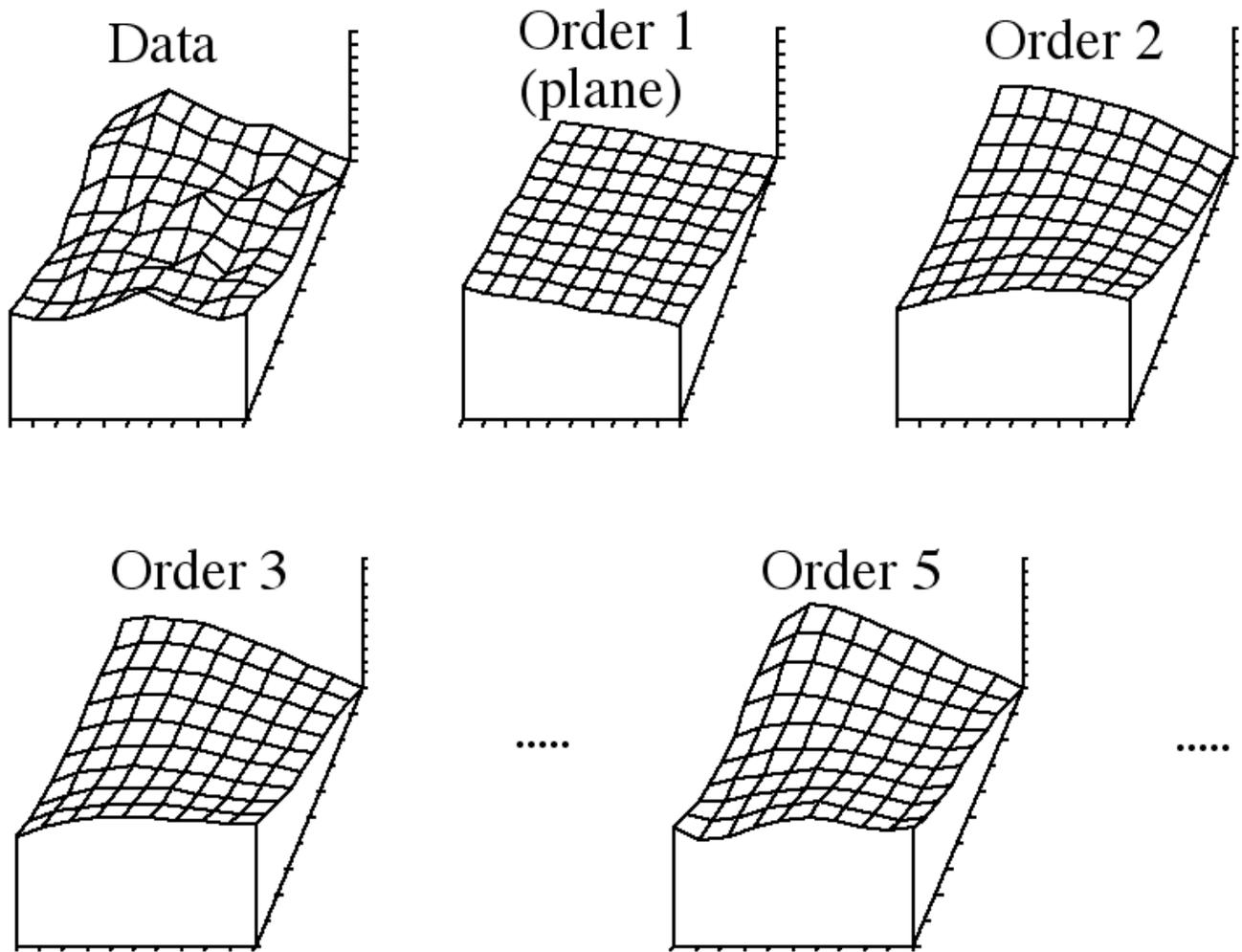


**Figure 45** - Density of an arthropod species along a double gradient and linear model.

If a plane does not explain enough variation, one can try to fit higher order polynomials, by adding second, third ...order  $x$ - $y$  terms and their products. The following equation is a cubic trend surface equation:

$$\hat{z} = b_0 + b_1x + b_2y + b_3x^2 + b_4xy + b_5y^2 + b_6x^3 + b_7x^2y + b_8xy^2 + b_9y^3$$

It is easy to visualize the outcome of the addition of one order to a trend surface model by remembering that each addition of an order allows one more fold to the surface (Figure 46):



**Figure 46** - Example of trend surface analysis, equations of order 1, 2, 3 and 5.

Trend surface analysis can model relatively simple structures with a reasonable amount of “hills” and “holes” resulting of one or two long-range trends (hence the name) across the sampling area. But this method, although easy to compute, suffers from several conceptual and practical problems, and should be used with great care. Here are some of these problems:

*Conceptual problem:*

- fitting a trend surface is useful only when the trend has an underlying physical or biological explanation, or if it can help

generating biological hypotheses; interpretation of individual terms is often difficult;

*Practical problems:*

- when data points are few, extreme values can seriously distort the surface;
- the surfaces are extremely susceptible to edge effects. Higher-order polynomials can turn abruptly near area edges, leading to unrealistic values;
- trend surfaces are inexact interpolators. Because they are long-range models, extreme values of distant data points can exert an unduly large influence, resulting in poor local estimates of the studied variable.

## **Detrending**

Despite its problems, trend surface analysis is very useful in one specific case. It has been said in Section 6.2.5 that, for testing, the condition of second-order stationarity or, at least, the intrinsic assumption must be satisfied. Removing a trend from the data at least makes the mean constant over the sampling area (although it does not address any problem of heterogeneity of variance). Furthermore, most methods of spatial analysis are devised to model the intermediate-scale component of spatial variation and are therefore much more powerful on detrended data. For these reasons, trend surface analysis is often used to detrend data: one fits a plane on the data and proceeds to analyze the finer scale structure on the residuals of this regression (this is equivalent to subtract the fitted values from the raw data and to work with what remains).