

# Nonlinear Regression Models and Applications in Agricultural Research

Sotirios V. Archontoulis and Fernando E. Miguez\*

## ABSTRACT

Nonlinear regression models are important tools because many crop and soil processes are better represented by nonlinear than linear models. Fitting nonlinear models is not a single-step procedure but an involved process that requires careful examination of each individual step. Depending on the objective and the application domain, different priorities are set when fitting nonlinear models; these include obtaining acceptable parameter estimates and a good model fit while meeting standard assumptions of statistical models. We propose steps in fitting nonlinear models as described by a flow diagram and discuss each step separately providing examples and updates on procedures used. The following steps are considered: (i) choose candidate models, (ii) set starting values, (iii) fit models, (iv) check convergence and parameter estimates, (v) find the “best” model among competing models, (vi) check model assumptions (residual analysis), and (vii) calculate statistical descriptors and confidence intervals. The associated feedback mechanisms are also addressed (i.e., model variance homogeneity). In particular, we emphasize the first step (choose candidate models) by providing an extensive library of nonlinear functions (77 equations with the associated parameter meanings) and examples of typical applications in agriculture. We hope that this contribution will clarify some of the difficulties and confusion with the task of using nonlinear models.

In **data analysis**, we often ask the following questions: Which is the best model to describe our data? Which is the best statistical index to judge the goodness of fit? How do we choose among competing models? There are no simple answers to these questions. Here we attempt to provide agronomists with a general framework on how to approach these questions appropriately. Our specific objectives are: (i) to provide a succinct overview of nonlinear models and to develop a guideline to understand the family of functions used in agricultural applications; (ii) to indicate techniques to modify nonlinear models and how to cope with multiple nonlinear models; (iii) to discuss key methodological issues on parameter estimation, model performance, and comparison; and (iv) to demonstrate step-by-step analysis of experimental data using a nonlinear regression model. The structure follows the flow diagram in Fig. 1. We start with the definition of nonlinear regression models and discuss their main advantages and disadvantages. Then we present 77 nonlinear functions (including those in supplemental tables) with references to applications in agriculture. We offer an updated overview of methodologies to fit models, choose starting values, assess goodness of fit, select the best models, and evaluate

residuals. Finally, we reanalyze experimental data on biomass growth with time (Danalatos et al., 2009).

## NONLINEAR REGRESSION MODELS

### Definition

In general, statistical models used in agricultural applications can be described with the following notation:

$$y = f(x, \theta) + \epsilon$$

where  $y$  is the response variable,  $f$  is the function or model,  $x$  are the inputs,  $\theta$  denotes the parameters to be estimated, and  $\epsilon$  is the error. Each parameter can be evaluated for whether it is linear or not: if the second derivative of the function with respect to a parameter is not equal to zero, then the parameter is nonlinear. Thus a given function ( $f$ ) can have a mix of linear and nonlinear parameters.

### Why Should We Use Nonlinear Models?

The main advantages of nonlinear models are parsimony, interpretability, and prediction (Bates and Watts, 2007). In general, nonlinear models are capable of accommodating a vast variety of mean functions, although each individual nonlinear model can be less flexible than linear models (i.e., polynomials) in terms of the variety of data they can describe; however, nonlinear models appropriate for a given application can be more parsimonious (i.e., there will be fewer parameters involved) and more easily interpretable. Interpretability comes from the fact that the parameters can be associated with a biologically meaningful process. For example, one of the most widely used nonlinear models is the logistic equation (Eq. [2.1] in Table 1). This model describes the pervasive S-shaped growth curve. The

Supplemental material available online. Dep. of Agronomy, Iowa State Univ., 1206 Agronomy Hall, Ames, IA 50011. Received 28 Dec. 2012. Accepted 18 Mar. 2013. \*Corresponding author (femiguez@iastate.edu).

Published in Agron. J. 107:786–798 (2015)

doi:10.2134/agronj2012.0506

Available freely online through the author-supported open access option.

Copyright © 2015 by the American Society of Agronomy, 5585 Guilford Road, Madison, WI 53711. All rights reserved. No part of this periodical may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or any information storage and retrieval system, without permission in writing from the publisher.

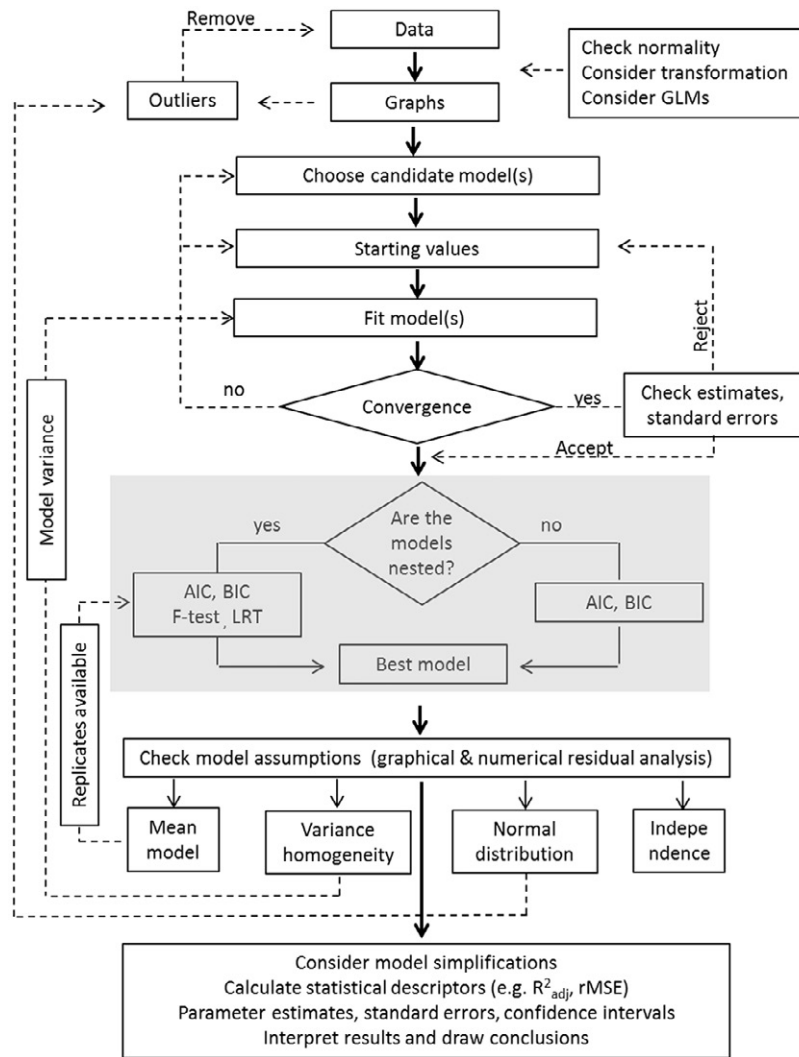


Fig. 1. Suggested work flow in the nonlinear regression analysis. Thick arrows indicate major steps, thin arrows indicate substeps, and dashed arrows indicate feedback in nonlinear regression. The shaded part is optional and can be ignored in simple cases. (Abbreviations: GLMs, generalized linear models; LRT, likelihood ratio test; AIC, Akaike information criterion; BIC, Bayesian information criterion.)

parameters have a clear meaning (see Table 1) and units associated with their definition. The asymptotic parameter ( $Y_{\text{asym}}$ ) has units equal to the response variable ( $Y$ ), the inflection point ( $t_m$ ) has units equal to the independent variable ( $t$ ), and the parameter that determines the steepness of the curves ( $k$ ) has units equal to  $t$ . This last parameter can be interpreted as the time (when  $t$  is time) that it takes to move from the inflection point to approximately 0.73 of the asymptotic value. A competing polynomial model used to describe the same data would have the disadvantages that more parameters would be needed (more than just three) and that the parameters would not be easily interpretable (Pinheiro and Bates, 2000). For example, what would be the interpretation of the parameters in a five-degree polynomial?

The final advantage of using nonlinear regression models is that their predictions tend to be more robust than competing polynomials, especially outside the range of observed data (i.e., extrapolation). Nonlinear regression models, however, come at a cost. Their main disadvantages are that they can be less flexible than competing linear models and that generally there is no analytical solution for estimating the parameters. The first point has as a consequence that the choice of model is crucial.

It is tempting to then try a large library of functions and choose the model with the lowest error; however, it is almost always better to choose a model based on whether it has been used successfully in similar applications and has biologically meaningful parameters (e.g., Table 1).

The lack of an analytical solution has two practical consequences. First, a numerical method needs to be used to find estimates for the parameters, and this implies that convergence of the algorithm needs to be checked (Fig. 1). A lack of convergence often results from the second consideration, which is that these numerical methods require starting values. Choosing a model with biologically meaningful parameters makes the process of choosing starting values easier because the starting values can usually be easily determined from visual inspection of the data (see below).

## TYPICAL NONLINEAR MODELS AND APPLICATION EXAMPLES

Choosing competing models for an application is not always a simple task. We have developed a reference table as a guideline to understand the family of functions used in agricultural applications. Table 1 presents 27 common nonlinear equations,

Table I. Nonlinear regression models. For example fits, see supplemental figures.

Eq.	Name and/or reference	Form	Parameter definition
Group I—Exponential			
[1.1]	Exponential decay	$Y = Y_o \exp(-kt)$	Y is the response variable (e.g., soil organic matter), t is the explanatory variable (e.g., time), $Y_o$ is the initial or the maximum Y value, k is a rate constant that determines the steepness of the curve
[1.2]	Exponential gives rise to maximum	$Y = Y_o[1 - \exp(-kt)]$	
Group II—Sigmoid functions			
[2.1]	Logistic (Verhulst, 1838)	$Y = Y_{\text{asym}}/[1 + \exp[-k(t - t_m)]]$	Y is the response variable (e.g., biomass), t is the explanatory variable (e.g., time), $Y_{\text{asym}}$ or $Y_{\text{max}}$ is the asymptotic or the maximum Y value, respectively, $t_m$ is the inflection point at which the growth rate is maximized, k controls the steepness of the curve, v deals with the asymmetric growth (if v = 1, then Richards' equation becomes logistic), a and b are parameters that determine the shape of the curve, $t_c$ is the time when $Y = Y_{\text{asym}}$ , $t_c$ is the critical time for a switch-off to occur (e.g., critical photoperiod), n is a parameter that determines the sharpness of the response
[2.2]†	Richards (1959)	$Y = Y_{\text{asym}}/[1 + v \exp[-k(t - t_m)]]^{1/v}$	
[2.3]	Gompertz (1825)	$Y = Y_{\text{asym}} \exp\{-\exp[-k(t - t_m)]\}$	
[2.4]	Weibull (1951)	$Y = Y_{\text{asym}}[1 - \exp(-at^b)]$	
[2.5]‡	Beta (Yin et al., 2003a)	$Y = Y_{\text{max}} \left( 1 + \frac{t_c - t}{t_c - t_m} \right)^{\left( \frac{t}{t_c} \right)^{Y_c/(t_c - t_m)}}$	
[2.6]§	Hill (switch-off) function	$Y = t_c^n/(t_c^n + t^n)$	
Group III—Photosynthesis			
[3.1]	Blackman (1905)	$Y = \min(Y_{\text{asym}}, al) - R_d$	Y is the response variable (net photosynthesis), l is the explanatory variable (irradiance), $Y_{\text{asym}}$ is the asymptotic Y value, a is the initial slope of the curve at low l levels, $R_d$ is the dark respiration, $\theta$ is a dimensionless curvature parameter (when $\theta = 1$ , Eq. [3.4] is equivalent to Eq. [3.1], and when $\theta \rightarrow 0$ , Eq. [3.4] is equivalent to Eq. [3.3])
[3.2]¶	Asymptotic exponential	$Y = Y_{\text{asym}}[1 - \exp(-al/Y_{\text{asym}})] - R_d$	
[3.3]¶	Rectangular hyperbola	$Y = alY_{\text{asym}}/(Y_{\text{asym}} + al) - R_d$	
[3.4]¶#	Nonrectangular hyperbola	$Y = \frac{Y_{\text{asym}} + al - \sqrt{(Y_{\text{asym}} + al)^2 - 4\theta alY_{\text{asym}}}}{2\theta} - R_d$	
[3.5]	Modified logistic (Sinclair and Horie, 1989)	$Y = Y_{\text{asym}}(2/[1 + \exp[-k(N - N_{\text{min}})]] - 1)$	Y is the response variable (light-saturated net photosynthesis), N is the explanatory variable (leaf N), $Y_{\text{asym}}$ is the asymptotic Y value, k determines the curvature of the curve, $N_{\text{min}}$ is the N value at or below which $Y = 0$
[3.6]††	Farquhar et al. (1980)	$Y = \min \left\{ \frac{V_{\text{cmax}}(C_i - \Gamma_*)}{C_i + K_{\text{mc}}(1 + O/K_{\text{mo}})}, \frac{J(C_i - \Gamma_*)}{4C_i + 8\Gamma_*} \right\} - R_{\text{day}}$	Y is the response variable (net photosynthesis), $C_i$ is the explanatory variable (intercellular $\text{CO}_2$ concentration), $V_{\text{cmax}}$ is the maximum carboxylation capacity, $\Gamma_*$ is the $\text{CO}_2$ compensation point in the absence of $R_d$ , $K_{\text{mc}}$ and $K_{\text{mo}}$ are Michaelis–Menten coefficients of Rubisco for $\text{CO}_2$ and $\text{O}_2$ , respectively, O is the partial pressure of $\text{O}_2$ (= 21 kPa), J is the photosystem II electron transport rate, $R_{\text{day}}$ is the dark respiration occurring in the light
Group IV—Temperature dependencies			
[4.1]	van't Hoff (1898) (known as the $Q_{10}$ function)	$Y = Q_{10}^{(T - T_{\text{ref}})/10}$	Y is the response variable (e.g. respiration), T is the explanatory variable (temperature), $T_{\text{ref}}$ is a reference temperature at which $Y = 1$ , $Q_{10}$ is the factor by which the rate of a process (respiration) increases for each 10°C temperature increase, E is the activation energy that determines the increase in temperature response, R is the universal gas constant (= 8.314 J K <sup>-1</sup> mol <sup>-1</sup> ), D is the deactivation energy that determines the decrease in the temperature response, S is the entropy term that determines the transition state of the curve, $E_o$ is an activation-energy-like parameter that is temperature adjusted, $T_x$ is a fitted temperature parameter (in K), $T_{\text{min}}$ is the base or minimum temperature for $Y = 0$
[4.2]	Arrhenius (1889)	$Y = \exp\{E/R[1/(T_{\text{ref}} + 273) - 1/(T + 273)]\}$	
[4.3]‡‡	Modified Arrhenius	$Y = \exp \left\{ \frac{E}{R} \left( \frac{1}{T_{\text{ref}} + 273} - \frac{1}{T + 273} \right) \right\} \times \left[ \frac{1 + \exp \{[(T_{\text{ref}} + 273)S - D/(T_{\text{ref}} + 273)R]\}}{1 + \exp \{(S/R) - [D/(T + 273)R]\}} \right]$	
[4.4]	Lloyd and Taylor (1994)	$Y = \exp\{E_o[l/(T_{\text{ref}} + 273 - T_x) - 1/(T + 273 - T_x)]\}$	
[4.5]	Ratkowsky et al. (1982)	$Y = (T - T_{\text{min}})^2/(T_{\text{ref}} - T_{\text{min}})^2$	
Group V—Peak or bell-shaped curves			
[5.1]	Beta (Yin et al. 1995)	$Y = \left[ \left( \frac{T_c - T}{T_c - T_o} \right) \left( \frac{T - T_b}{T_o - T_b} \right) \right]^{(T_o - T_b)/(T_c - T_o) \cdot c}$	Y is the response variable (e.g., rate of development), T is the explanatory variable (temperature), $T_o$ is the optimum temperature for maximum Y, $T_b$ is the base or minimum temperature for $Y = 0$ , $T_c$ is the ceiling or maximum temperature for $Y = 0$ , c is a curvature parameter (default c = 1)
[5.2]§§	Bell curve	$Y = Y_{\text{asym}} \exp[a(X - X_o)^2 + b(X - X_o)^3]$	Y is the response variable, X is the explanatory variable, $Y_{\text{asym}}$ is the asymptotic maximum Y value, $X_o$ is the position of the center of the peak ( $Y_{\text{asym}}$ ), a (default = 0.5 for the Gaussian function), and b are coefficients controlling the width of the bell
[5.3]	Gaussian function	$Y = Y_{\text{asym}} \exp\{-0.5[(X - X_o)/b]^2\}$	

(continued)

Table 1. Continued.

Eq.	Name and/or reference	Form	Parameter definition
Group VI—Other nonlinear equations			
[6.1]	Power	$Y = aX^b$	Y is the response variable, X is the explanatory variable, a and b are parameters that define the shape of the curve and the magnitude of the Y value
[6.2]	Modified hyperbola	$Y = aX/(1 + bX)$	
[6.3]	Michaelis–Menten	$Y = \mu_{\max}X/(X + C_{\text{sat}})$	Y is the response variable (e.g., denitrification rate), X is the explanatory variable—the substrate (e.g., $\text{NO}_3^-$ ), $\mu_{\max}$ is the rate constant, $C_{\text{sat}}$ is the half-saturation constant
[6.4]¶¶¶	Rational function	$Y = a_1X^{a_2}/(1 + a_3X^{a_4})$	Y is the response variable, X is the explanatory variable, $a_1$ and $a_3$ are parameters defining the magnitude of the Y value, $a_2$ and $a_4$ are parameters defining the shape of the curve (if $a_2 = a_4 + 1$ , then the equation shows a near-linear response; if $a_2 = a_4$ , then the equation becomes a hyperbola; if $a_2 < a_4$ , then the equation takes a bell shape; if $a_2 > a_4$ or $a_4 = 0$ , then the equation becomes exponential; if $a_2 = 0$ , then the equation becomes exponential decay
[6.5]	Ricker curve	$Y = a_1X \exp(-a_2X)$	Y is the response variable, X is the explanatory variable, $a_1$ and $a_2$ are parameters that control both the height and the width of the right skew of the “bell”

† The maximum growth rate for the Richards equation is given in Birch (1999).

‡ The maximum growth rate for the Beta equation is given in Yin et al. (2003a).

§ Cited in Amaducci et al. (2008).

¶ Cited in Goudriaan (1979).

# Goudriaan (1979) and Johnson et al. (2010) used the nonrectangular hyperbola to describe the photosynthetic rate response to  $\text{CO}_2$ .

†† This is simplified version of the Farquhar model.

‡‡ The optimum temperature for this equation is given in Medlyn et al. (2002).

§§ Cited in Hammer et al. (2009).

¶¶¶ Cited in Bril et al. (1994).

and Supplementary Tables S2 to S6 supplement these with 45 additional equations. We classified the equations into six groups based on a combination of statistical form and use in the agricultural domain. All equations have been used in agricultural applications, and most of the parameters have an interpretable meaning (see supplementary figures also). The variety of equations presented in Table 1 reflects well the fact that one equation does not suit all processes.

### Group I—Exponential

The exponential decay and exponential gives rise to maximum functions (Eq. [1.1] and [1.2], Table 1) find applications in a wide spectrum of soil and plant sciences. They are commonly used to describe light and N vertical distributions within plant canopies (Monsi and Saeki, 2005),  $\text{N}_2\text{O}$  emission response to N fertilizer (e.g., Hoben et al., 2011), cumulative soil respiration (e.g., Gillis and Price 2011), photoperiodic sensitivity (e.g., Wang and Engel, 1998), temperature or moisture responses to nitrification (e.g., Ma and Shaffer, 2001), water infiltration rate (Horton, 1940), and first-order kinetics. They are simple equations with one major unknown, the rate constant ( $k$ ), which is also termed the *extinction coefficient* in crop physiology. The ratio  $\ln(2)/k$  is of importance in soil science because it denotes the mean residence time (e.g., soil organic matter). Equation [1.1] provided the starting point to develop case-specific nonlinear functions. Yin et al. (2000, 2003b) established a nonlinear function to describe leaf area index development as a function of canopy N content (Eq. [1.8] in Supplementary Table S1). Johnson et al. (2010) developed a flexible nonlinear function for the protein (N) vertical distribution within plant canopies (Eq. [1.9] in Supplementary Table S1). In soil science, Andren and Paustian (1987) and Gillis and Price (2011) extended Eq. [1.1] to better describe the

decomposition of straw residue and biochar, respectively (see Supplementary Table S1). Lastly, Eq. [1.1] as well as exponential functions (viz. Eq. [2.16] and [2.15] in Supplementary Table S2) were also applied to describe the initial parts of growth curves but not the entire growth curves because the growth profile often reaches an asymptotic value.

### Group II—Sigmoid Curves

Sigmoid curves (mathematical functions having an S shape) are another important group of nonlinear models. These models are often applied to describe plant height, weight, leaf area index, or seed germination as a function of time, N application rate, herbicide dose, etc. (e.g., Gan et al., 1996; Miguez et al., 2008). Sigmoid equations are also used as 0–1 modifiers in process-based models to incorporate moisture availability or soil pH, etc., effects on soil N transformation processes (e.g., McGechan and Wu, 2001) and also as a switch-off function in studies assessing plant photoperiodic sensitivity (e.g., Amaducci et al., 2008). Table 1 presents common sigmoid functions and Supplementary Table S2 provides additional sigmoidal equations, providing increased flexibility (e.g., when maximum growth or the inflection point is achieved at the start or the end of growth period). Additional equations can be found in Zwietering et al. (1990), Zeide (1993), Leduc and Goelz (2009), and many statistical textbooks or software manuals (e.g., SigmaPlot, JMP, TableCurve).

In general, the suitability of a sigmoid equation to estimate maximum rate of increase or optimum x level for maximizing the y value is an important part of its function (Birch, 1999; Yin et al., 2003a). Each function has its advantages and disadvantages (for a discussion, see Birch, 1999; Yin et al., 2003a), and it is up to the researcher to select the most appropriate one to fit the experimental data. The logistic equation, Eq. [2.1], describes symmetric growth having an inflection point at half



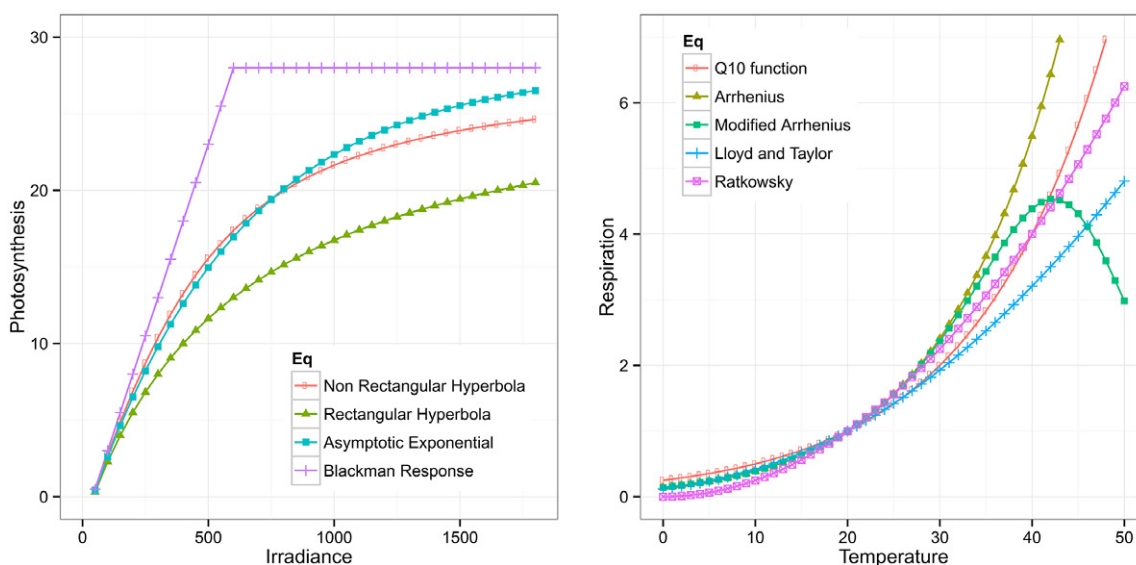


Fig. 2. Nonlinear models for describing photosynthesis response to irradiance (left) and respiration response to temperature (right). Equations are given in Table 1. The following parameter values were used for these plots: asymptotic maximum response variable ( $Y_{\text{asym}}$ ) = 30  $\mu\text{mol CO}_2 \text{ m}^{-2} \text{ s}^{-1}$ , initial curve slope ( $a$ ) = 0.05  $\text{mol CO}_2 \text{ mol}^{-1} \text{ photons}$ , curvature parameter ( $\theta$ ) = 0.7, dark respiration ( $R_d$ ) = 2  $\mu\text{mol CO}_2 \text{ m}^{-2} \text{ s}^{-1}$ ; increase in respiration for each 10°C temperature increase ( $Q_{10}$ ) = 2, reference temperature ( $T_{\text{ref}}$ ) = 20°C, universal gas constant ( $R$ ) = 8.314  $\text{J K}^{-1} \text{ mol}^{-1}$ , activation energy ( $E$ ) = 65,000  $\text{J mol}^{-1}$ , entropy ( $S$ ) = 650  $\text{J K}^{-1} \text{ mol}^{-1}$ , deactivation energy ( $D$ ) = 207,000  $\text{J mol}^{-1}$ , temperature-adjusted activation-energy-like parameter ( $E_0$ ) = 350 K, fitted temperature parameter ( $T_x$ ) = 225 K, and minimum temperature ( $T_{\text{min}}$ ) = 0°C. Note that at the reference temperature of 20°C, respiration = 1. The optimum temperature for the modified Arrhenius equation is:  $T_{\text{opt}} = D/[S - R \ln[E/(D - E)]] - 273 = 42.3^\circ\text{C}$ .

the final size. The Gompertz equation, Eq. [2.3], has an inflection point that is controlled by its asymptotic value and is at about one-third ( $1/e = 0.3679$ ), while others like the Richards or Weibull or beta have more flexibility in dealing with asymmetric growth (the inflection point can be at any  $x$  value).

Having a flexible inflection point is another important feature of a sigmoid curve. For that, Birch (1999), for example, modified the logistic equation (Eq. [2.1]) to deal with asymmetric growth by adding an extra shape parameter. When growth is known to decrease after a certain period of time, then the beta function (Eq. [2.5]) might be a better option (see supplementary figures). On the other hand, Eq. [2.5] might not accurately predict initial growth and, in cases when the initial phase is very important, different versions of the beta function should be used (see Eq. [2.11] in Supplementary Table S2 and example below). It is important to note that all sigmoid equations presented in Table 1 (except Eq. [2.4] and [2.5]) assume an initial  $Y$  value close to zero at time zero, which is reasonable in most cases, e.g., at planting, the biomass weight is very close to zero.

### Group III—Photosynthesis

Photosynthesis is the most important biological process involved in plant growth, and its rate is influenced by irradiance, temperature, N availability, the vapor pressure deficit, and  $\text{CO}_2$  concentration. Different nonlinear functions have been developed to describe the photosynthesis response to different environmental variables. Functions to describe the photosynthesis response to irradiance have been researched the most (Jassby and Platt, 1976; Goudriaan, 1979). All equations assume that dark respiration ( $R_d$ ) is independent of the light level. Among the equations presented in Table 1, Blackman (Eq. [3.1]) is the simplest one, and the asymptotic exponential (Eq. [3.2]) and the nonrectangular hyperbola (Eq. [3.4]) are the most common. The rectangular hyperbola (Eq. [3.3], also termed the Michaelis–Menten equation) is used less frequently because it reaches saturation faster than photosynthesis actually does.

Currently, the scientific discussion on the photosynthetic capacity ( $Y_{\text{asym}}$ ) and efficiency ( $a$ ) of different plant species is based on the comparison of nonlinear regression estimates; for this reason, caution should be exercised because similar estimates from different equations can result in different responses (Fig. 2). Equation [3.2] is a simple three-parameter equation widely used in light-driven process-based models like SUCROS and Hybrid-maize (Goudriaan and van Laar, 1994; Yang et al., 2004). Equation [3.4] offers more flexibility and is more accurate than Eq. [3.2] at the cost of one extra parameter (i.e.,  $\theta$ , the curvature parameter). When  $\theta = 1$ , Eq. [3.4] becomes the Blackman equation (Eq. [3.1]), and when  $\theta$  approaches zero, Eq. [3.4] becomes the rectangular hyperbola equation (Eq. [3.3]). Equation [3.4] is the reference equation when the biochemical model of Farquhar et al. (1980) or Collatz et al. (1992) is used in modeling studies. New equations are still being developed and tested (e.g., Eq. [3.7] in Supplemental Table S3).

The photosynthesis response to  $\text{CO}_2$  has been quantified empirically using a nonrectangular hyperbola (Goudriaan, 1979; Johnson et al., 2010) and mechanistically using a biochemical model (Farquhar et al., 1980). The biochemical model is based on Michaelis–Menten kinetics for substrate-limited growth and the law of minimum between carboxylation and electron transport rates (Eq. [3.6], Table 1). Although its computation is laborious, this model has found large acceptance. For more details on that model, see the original publications (Farquhar et al., 1980; von Caemmerer and Farquhar, 1981) and model application studies (Medlyn et al., 2002; Archontoulis et al., 2012). The photosynthesis response to leaf N, which is strongly related to the Rubisco content, can be modeled using a modified logistic equation proposed by Sinclair and Horie (see Eq. [3.5], Table 1), while alternatives exist (Eq. [3.8] in Supplemental Table S3). The photosynthesis response to water stress is usually described by sigmoid functions at the leaf level. For instance, Vico and Porporato (2008) utilized a

Weibull-type curve (Eq. [3.9] in Supplemental Table S3). The photosynthesis response to a vapor pressure deficit has been described by an exponential decay function (e.g., Osório et al., 2006) but usually more sophisticated approaches have been used (Collatz et al., 1992; Yin and Struik, 2009). The photosynthesis response to temperature is discussed below.

### Group IV—Temperature Dependence

A multitude of nonlinear regression models have been proposed and tested for modeling the temperature dependence of various soil and plant processes (Lloyd and Taylor, 1994; Kätterer et al., 1998; Davidson et al., 2006; Shibu et al., 2006; Portner et al., 2010). These include power, logarithmic, exponential, sigmoid, and bell-shape functions (Table 1; Supplemental Table S4). The van't Hoff or  $Q_{10}$  function ( $Q_{10}$  is the factor by which the rate of a process increases for each 10°C temperature increase) has found application in many studies, particularly in those addressing leaf or soil respiration rates. A  $Q_{10}$  of 1 indicates no temperature effect. The  $Q_{10}$  value commonly ranges from 1.4 to 4.9 (Tjoelker et al., 2001; Atkin et al., 2005). In the Arrhenius equation, the  $Q_{10}$  term has been replaced by the activation energy. Both equations are equivalent, producing similar temperature responses (Fig. 2); however, it should be noted that both  $Q_{10}$  and  $E$  coefficients are temperature-range dependent. Usually, narrow temperature measurement ranges result in high and sometimes unrealistic  $Q_{10}$  or  $E$  estimates. Lloyd and Taylor (1994) noticed limitations of these two functions (i.e., the rate of reaction is not constant across temperatures) and developed a new equation (see Eq. [4.4]) to fit extensive literature data.

The above temperature functions describe a monotonic increase (Fig. 2). The rate of a process probably increases to an optimum temperature point and then drops (in reality, due to the lack of appropriate data, the drop is not always apparent). New equations or modifications of existing models have been developed to account for this. For example, the modified Arrhenius function, when compared with the Arrhenius equation, includes an additional two-parameter term (see  $D$  and  $S$  in Eq. [4.3] and Fig. 2) to capture the decline in the rate of a process at very high temperature (e.g., the electron transport rate). If one of the two additional parameters is set to zero, then Eq. [4.3] becomes Eq. [4.2] (see also supplemental figures). This equation is “fragile” and requires careful parameterization (Medlyn et al., 2002; Archontoulis et al., 2012).

Johnson et al. (2010) argued that temperature functions based on the activation energy of chemical reactions are quite complex and difficult to apply routinely. They used a modified beta function to describe the photosynthesis response to temperature (Eq. [4.10] in Supplemental Table S4). Kirschbaum (1995) used a modified exponential temperature function (see Eq. [4.9] in Supplemental Table S4) that provides a peak pattern to fit soil organic matter decomposition data. Additional (but difficult to interpret) peak temperature response functions were reported in Portner et al. (2010).

### Group V—Bell Curves

In addition to temperature dependencies of photosynthesis, the bell-shaped or peak functions have been applied in agricultural science to describe the rate of phenological development as a function of temperature (e.g., Yin et al., 1995), the size of

a leaf as a function of its rank in a plant (e.g., Hammer et al., 2009) or soil moisture effects on  $N_2O$  emissions (e.g., Rafigue 2011). Table 1 lists three important equations. More application examples and different types of bell-curve equations can be found in Ma and Shaffer (2001) and in Supplementary Table S5. In process-based simulation models, researchers have approximated a bell-shaped response (viz. rate of development) with two-, three-, or four-segment (broken) linear regression models (e.g., APSIM; Keating et al., 2003). Typically, these segmented models should be fit using nonlinear methods as well.

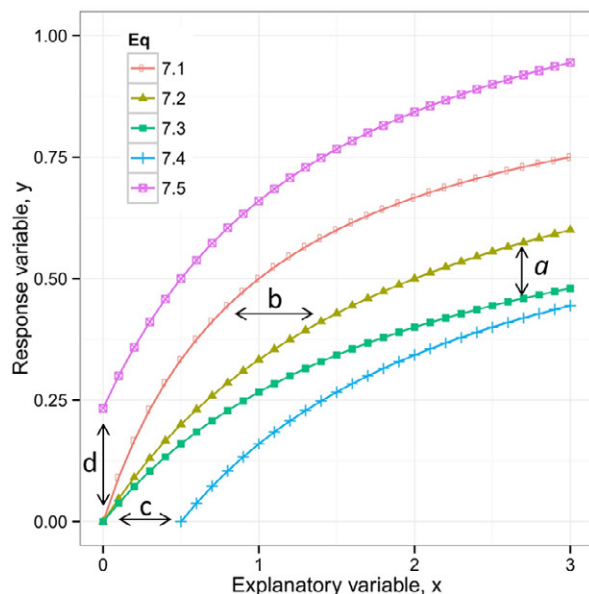
### Group VI—Others

In allometric studies, the relations that exist among the growth rates of different plant components are quantified by means of regression analysis. Given the large variability that exists among plant species and plant components, numerous nonlinear models have been utilized including power (Eq. [6.1], e.g., plant N concentration vs. biomass weight), hyperbolic (Eq. [6.2]), and sigmoid curves (e.g., Eq. [3.1]). For application examples, see Vega et al. (2000), Vega and Sadras (2003), and Archontoulis et al. (2010). The Michaelis–Menten equation (Eq. [6.3]) is well known and routinely applied to quantify the rate of a process (i.e., denitrification) that is dependent on the substrate (i.e.,  $NO_3$ ). In contrast, Eq. [6.4] is not as common in agronomy, but it appears to be very flexible, taking many forms from linear to exponential and bell curved (see supplemental figures). It was applied to model temperature effects on soil N mineralization (Bril et al., 1994). The last equation in Table 1 is the Ricker function (Eq. [6.5]), an option for hump-shaped patterns that are skewed to the right (Bolker, 2008).

### Manipulating or Combining Nonlinear Functions

Sometimes there is a need to modify a “standard” nonlinear function to fit a set of data. This has led to the development of numerous versions of a standard equation (e.g., Birch, 1999; Tsoularis, 2001; Supplemental Tables S1–S3). Using the simplest form of the Michaelis–Menten hyperbolic function (see Eq. [7.1] in Fig. 3), we illustrate simple modification techniques. Equation [7.1] starts at zero when  $x = 0$  and increases up to an asymptotic value of 1 as  $x$  increases. We can change the horizontal scale of this function by multiplying the variable  $x$  by a constant parameter,  $b$ , which is called a *scale parameter* (Bolker, 2008). If  $b > 1$ , then  $y$  saturates faster and if  $0 < b < 1$ , then  $y$  saturates more slowly (Eq. [7.2] in Fig. 3). We can change the vertical scale of the function by introducing a new parameter,  $a$  (Eq. [7.3] in Fig. 3). In this case, the asymptote moves from 1 to  $a$ . We can shift the whole curve to the right or the left by subtracting or adding a new parameter,  $c$ , to the  $x$  variable (Eq. [7.4] in Fig. 3), which is called the *location parameter* (Bolker 2008). Similarly we can shift the whole curve upward or downward by adding or subtracting a new constant value,  $d$  (Eq. [7.5] in Fig. 3). Lastly, we can replace  $x$  with  $x^k$ , where  $k$  is a shape parameter, and then the equation takes many forms (exponential, sigmoid, etc.; not shown). A close example to the last modification is Eq. [2. 6] in Table 1. When we modify nonlinear functions, we should add parameters that have an interpretable meaning.

When nonlinear functions are extended or combined to describe a phenomenon, we should be aware that there is an



(#)	Equation	Parameters
7.1	$Y = \frac{x}{1+x}$	-
7.2	$Y = \frac{bx}{1+bx}$	$b=0.5$
7.3	$Y = a \frac{bx}{1+bx}$	$b=0.5, a=0.8$
7.4	$Y = a \frac{b(x-c)}{1+b(x-c)}$	$b=0.5, a=0.8, c=0.5$
7.5	$Y = a \frac{b(x-c)}{1+b(x-c)} + d$	$b=0.5, a=0.8, c=0.5, d=0.5$

Fig. 3. Example of a nonlinear model modification. Starting with Eq. [7.1], the parameters  $a$ ,  $b$ ,  $c$ , and  $d$  were added step by step to Eq. [7.1], resulting in four new equations: Eq. [7.2–7.5]. Horizontal or vertical arrows in the figure panel indicate how the additional parameters affected the model.

upper limit in the number of parameters that can be estimated from standard nonlinear regression analysis. This depends on the complexity of the model and the number of data points. For example, to fit growth curves with three parameters, we need at least four data points. When a process is described by a combination of nonlinear models (e.g., Farquhar model of photosynthesis or generic simulation crop models), then a stepwise parameterization method is usually applied. For application examples, see Miguez et al. (2009) and Archontoulis et al. (2012).

## FITTING NONLINEAR MODELS

Presently there are many statistical software packages available for fitting nonlinear models (e.g., SAS, R, JMP, GenStat, MatLab, Sigmaplot, OriginLab, and SPSS). Nonlinear parameter estimates can be obtained using different methods (Bates and Watts, 2007); the most common are: (i) ordinary least squares, which minimizes the sum of squared error between observations and predictions, and (ii) the maximum likelihood method, which seeks the probability distribution that makes the observed data most likely. For non-normal data such as binomial or counts, generalized (non)linear models should be used (Lindsey, 2001; Huet et al., 2003; Gbur et al., 2012). Most problems encountered during the use of standard nonlinear regression software functions are due to a poor choice of competing models or an incorrect equation or starting values (Fig. 1). The choice of estimation method can affect the parameter estimates (Ruppert et al., 1989), but in general, estimates from least squares and maximum likelihood methods tend to differ only when the data are not normally distributed and are approximately identical when the data follow a normal distribution (Myung, 2003).

### Choosing Starting Values

All the procedures for nonlinear parameter estimation require initial values. The choice of values will influence the convergence of the estimation algorithm, in the worst case yielding no convergence and in the best case convergence in a few iterations (Ritz

and Streibig, 2005); however, there is no standard procedure for getting initial estimates. We indicate five practical methods:

1. If the model has parameters with biological meaning, then use information from the literature.
2. Use graphical exploration (see the example below and Fig. 4).
3. Transform the nonlinear model into a linear model. For instance logarithmic transformation of Eq. [1.1] yields a linear equation (viz.  $\ln Y = Y_0 - kt$ ) in which rough estimates of the parameter values can be easily obtained by linear regression. This method is recommended for getting initial estimates and to detect deviations from linearity, but these estimates may also be used as the final estimates (Ruppert et al., 1989). For more transformation examples, see Zeide (1993), Singh (2006), and Portner et al. (2010).
4. In the case where no clear guidelines exist for choosing starting values, the recommendation is to use a grid search or “brute force” approach (e.g., PROC NLIN in SAS or the nls2 package in R). This grid search can be done by generating an extensive coverage of possible parameter values (and their combinations) and then evaluating the model at each one of these parameter combinations. The numerical method can then be used starting with the combination that resulted in the best fit (lowest mean squared error). The hope is that an extensive enough coverage of the parameter space will provide a combination of parameters that will result in an adequate fit.
5. Use prespecified algorithms. This approach is specific to a given equation and can be used to calculate starting values for a given data set (e.g., Pinheiro and Bates, 2000; Ritz and Streibig, 2008).

### Checking Algorithm Convergence

After the initial attempt at fitting a nonlinear model, we recommend that algorithm convergence is evaluated (Fig. 1). Convergence is achieved when a measure (such as the relative offset or maximum change among parameter estimates; Bates and Watts, 2007) is below a certain threshold value (e.g.,  $10^{-5}$ ), meaning

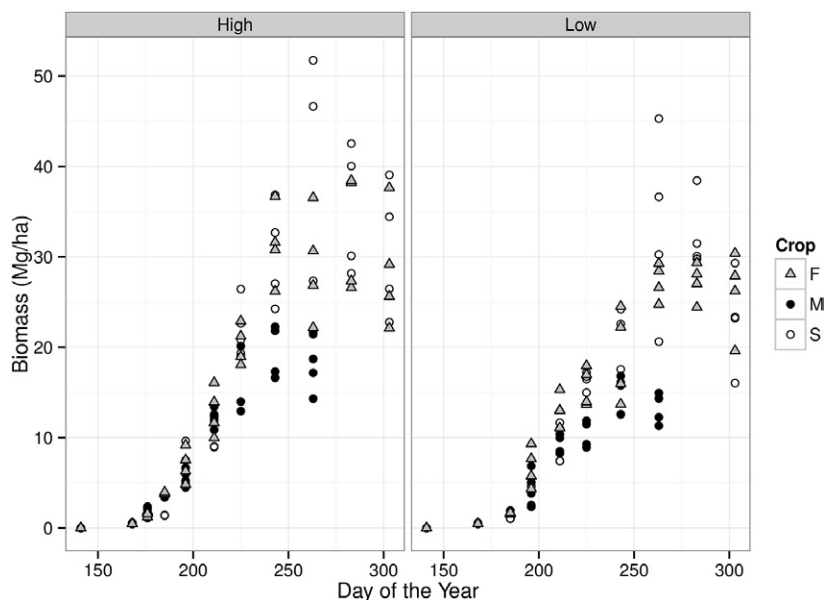


Fig. 4. Biomass accumulation with time for three crops—maize (M), fiber sorghum (F), and sweet sorghum (S)—at high and low levels of agricultural inputs, collected in Greece in 2008.

that the algorithm has found a “best” solution (Fig. 1). If convergence is not achieved, the most likely problems are a poor choice of starting values or the selected model is not well suited to describe the data. If convergence is achieved, the next step is to evaluate whether the parameter estimates are within a reasonable range. This requires not only evaluating the point estimates but also their standard errors. Unusually large standard errors are a sign of convergence problems, even if convergence was apparently achieved in the previous step. If no problems were encountered up to this point, the analysis can continue by assessing model assumptions and simplifying the model.

### Evaluating Model Assumptions

When we are dealing with one model, the next step is to evaluate key model assumptions: normally distributed errors, independent errors, and homogeneous variance for the errors (Fig. 1). This step and the following steps are not unique to nonlinear models but are common to all linear models. Substantial deviations from the assumptions could result in bias (inaccurate estimates), distorted standard errors, or both (Ritz and Streibig, 2008). Violations of these assumptions can be detected from an analysis of the residuals by means of graphical procedures and formal statistical tests. For a thorough analysis, see Ritz and Streibig (2008).

Briefly, to check whether the distribution of the measurement errors follows normality, the standardized residual plot is commonly applied (Pinheiro and Bates, 2000; see also the example later and Fig. 5). Outliers and many extreme values are common causes for deviations from normality (Fig. 1). Heterogeneity of variance can be detected by looking at the plot of the fitted values over the residuals (absolute residuals, which are raw residuals stripped of the negative sign, or standardized residuals, which are raw residuals scaled by the variance; see the example below).

When the residual errors show a trend (e.g., increasing variability as the explanatory variable increases, Fig. 4 and 5), this can be addressed by modeling the variance as a function of the

independent variable or the fitted values (Fig. 1 and 6). This is the case in our example (see discussion below). If variance heterogeneity is ignored, the parameter estimates might not be influenced much, but this may result in severely misleading confidence and prediction intervals (Carroll and Ruppert, 1988). The residuals are assumed to be independent, and when this assumption is violated it is visually evident in a plot of correlations of residuals against “lag” (or units of separation in time or space). Typically, variables measured with time on the same subject (e.g., plant, animal, or soil sample) tend to result in autocorrelated residuals that need to be accounted for by modeling the variance–covariance matrix.

### MODEL SELECTION CRITERIA

When we are dealing with multiple models, the question is how to find the best model among competing models. Depending on the structure of the models, different statistical criteria can be used to find the best model: *F* test, Akaike information criterion (AIC), Bayesian information criterion (BIC), or the likelihood ratio test (Zucchini, 2000; Burnham and Anderson, 2002; Hoffmann, 2005; Ritz and Streibig, 2008; Lewis et al., 2011). When models are *nested* (one model is a special case of another), any of these criteria are applicable (Fig. 1). When models are *non-nested* (models having different structures, e.g., Eq. [2.1] vs. Eq. [2.2]), typically the AIC and the BIC criteria are used (Fig. 1). From a practical point of view, however, one model might be preferred over another based on interpretability and specific objectives. There needs to be a balance between statistical model performance and how effectively the model answers research questions.

For two nested models, one with two parameters (reduced, e.g., Eq. [2.9] in Supplementary Table S2) and one with four parameters (full, e.g., Eq. [2.7] in Supplementary Table S2), to check whether the addition of parameters has a statistically significant contribution to the model performance, we can use the *F* test:



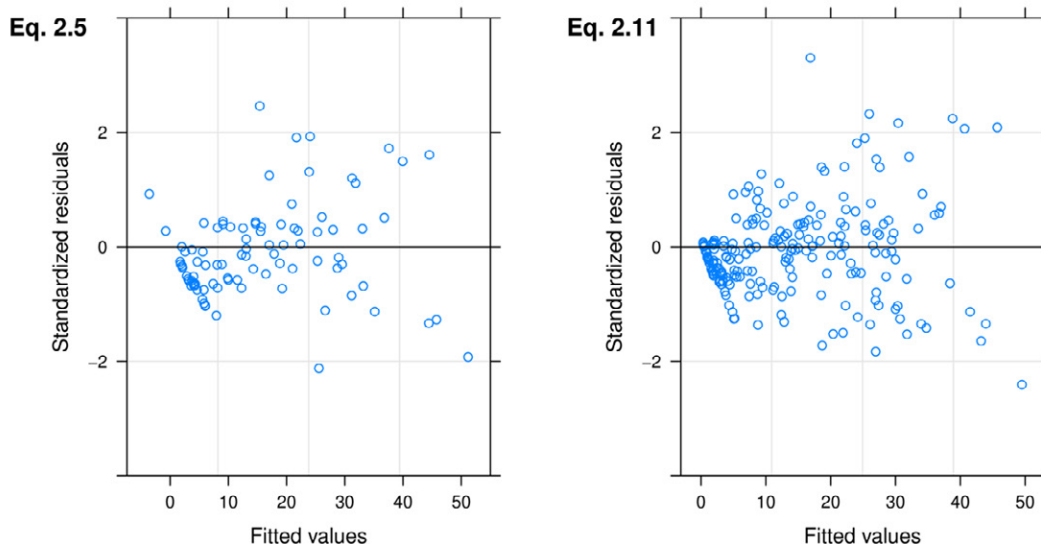


Fig. 5. Standardized residuals from individually fits to all experimental units: Eq. [2.5] from Table 1 (left) and Eq. [2.11] from Supplementary Table 2 (right). The fewer points in the left panel are because Eq. [2.5] converged in only 10 out of the 24 experimental units.

$$F = \frac{(SS_{\text{full}} - SS_{\text{reduced}}) / (df_{\text{full}} - df_{\text{reduced}})}{SS_{\text{full}} / df_{\text{full}}}$$

where  $SS_{\text{full}}$  and  $SS_{\text{reduced}}$  are the regression model sum of squares for the full and reduced models, respectively, and  $df_{\text{full}}$  and  $df_{\text{reduced}}$  are the degrees of freedom for the full and the reduced models, respectively. The  $P$  value can be calculated at  $df_{\text{full}} - df_{\text{reduced}}$ ,  $n - p - 1$ , (equals 2,  $n - 3$  for this example), where  $p$  is the number of parameters for the full model and  $n$  is the number of observations, and a decision can be made. This test is sometimes referred to as extra-sum-of-squares or multiple partial  $F$  test (Hoffmann, 2005; Ritz and Streibig, 2008). The  $F$  test is computed when the ordinary least squares method is used to fit the data (see fitting nonlinear models above). When the maximum likelihood method is used to fit the data, then the likelihood ratio test statistic ( $Q$ ) is computed to compare nested models:

$$Q = 2 \log \left( \frac{L_{\text{full}}}{L_{\text{reduced}}} \right) \\ = 2 (\log L_{\text{full}} - \log L_{\text{reduced}})$$

where  $L_{\text{full}}$  and  $L_{\text{reduced}}$  are the likelihood functions for the full and reduced models, respectively. These functions are closely related to the residual sum of squares (see Eq. [2.4] in Ritz and Streibig, 2008). It is assumed that  $Q$  is approximately  $\chi^2$  distributed with  $n - p - df_{\text{reduced}}$  degrees of freedom (for details, see Ritz and Streibig, 2008). Another approach for model selection involves calculating the AIC and BIC values for each model separately:

$$AIC_i = -2 \log L_i + 2p_i$$

$$BIC_i = -2 \log L_i + p_i \log n$$

where  $L_i$  and  $p_i$  are the likelihood and the number of parameters for each model, and  $n$  is the number of observations. For both statistical criteria, a smaller value indicates a preferable model. The BIC differs from the AIC only in the second term, which depends on  $n$ . Clearly as  $n$  increases, the BIC favors the simpler models (fewer parameters). This explains why sometimes the AIC and BIC indices disagree. For more information about these indices, see Burnham and Anderson (2002). Note that the likelihood ratio test, AIC, and BIC are all designed to compare the performance of models that have been fitted to data via maximum likelihood estimation (or for any model for which the likelihood can be calculated).

### Goodness of Fit

There is no single method or index to best assess the goodness of fit, but there are many different methods (graphical and numerical) that highlight different features of the data and the model. Graphical comparison provides a quick visual assessment of the goodness of fit. Numerical statistical indices like  $R^2$ , adjusted  $R^2$  ( $R^2_{\text{adj}}$ ), bias, mean squared error, root mean squared error (RMSE), modeling efficiency (ME), concordance correlation, and others (Wallach, 2006) provide the additional detail needed to assess the goodness of fit. Some indices measure the absolute error (includes units) and some others the relative error (excludes units). Depending on the data type, a combination of these indices can be used. For example, the relative term is more meaningful than the absolute when comparing errors based on different data sets. An important aspect of the statistical descriptors is that some simple and very common indices like  $r^2$  and bias do not account for the number of parameters. The following numerical indices are commonly used in model evaluation:

$$\text{bias} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)$$

$$R^2 = 1 - \frac{SS_{\text{residual}}}{SS_{\text{total}}}$$

$$= 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}$$

$$R^2_{\text{adj}} = \left( R^2 - \frac{p}{n-1} \right) \left( \frac{n-1}{n-p-1} \right)$$

$$\text{RMSE} = \sqrt{\frac{SS_{\text{residual}}}{n-p-1}}$$

$$\text{ME} = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

where bias,  $R^2$ ,  $R^2_{\text{adj}}$ , RMSE, and ME are numerical statistical indices,  $n$  is the number of data points,  $Y_i$  and  $\hat{Y}_i$  are the observed and predicted values, respectively,  $\bar{Y}$  is the mean observed value,  $p$  is the number of model parameters, and  $SS_{\text{residual}}$ ,  $SS_{\text{total}}$  are the sum of squares for the residual, regression model, and total, respectively.

Although often used, the  $R^2$  does not represent a good metric of model performance for nonlinear models. It has several limitations (e.g., it does not account for the number of parameters), and other measures of agreement (or combinations) should be used (Wallach, 2006). The main limitation of  $R^2$  is that the full model does not necessarily include the simpler model with one single parameter, as is the case with linear models.

The numerical statistical descriptors indicate the average performance of the model across the sample. When the variability is not constant throughout the sample (e.g., biomass increase with time), then statistical indices do not capture the fact that the uncertainty is not the same at different magnitudes of the response variables (see the example below; Fig. 4 and 6). We should often be concerned with the predictive ability of the model, and for this, cross-validation techniques can be used and the mean squared error of prediction is more appropriate (Wallach, 2006).

## EXAMPLE APPLICATION

The example follows the workflow illustrated in Fig. 1.

**Data:** We used data from Danalatos et al. (2009), which represent destructive measurements of aboveground biomass accumulation with time for three crops: fiber sorghum (F), sweet sorghum (S), and maize (M), growing in a deep fertile loamy soil of central Greece under two management practices: high and low input conditions, in 2008. High refers to weekly irrigation (to match 100% of maximum evapotranspiration) and application of 200 kg N ha<sup>-1</sup> and low input refers to biweekly irrigation (approximately 50% of maximum evapotranspiration) and application of 50 kg N ha<sup>-1</sup>. The experiment was a 2 × 3 factorial completely randomized in four blocks. For more details, see Danalatos et al. (2009). With such data, many questions are possible. We will concentrate on three: (i) what is the maximum biomass accumulated by these crops, (ii) at what point in time was this biomass achieved, and (iii) are there significant treatment effects and/or interactions. This requires statistical determination of the effects of crop type on the function parameters and also the effect of input level (i.e., high or low). These questions are approachable through the use of a nonlinear model that captures the mean function and the structure of the data.

**Graphs:** Visually (Fig. 4), sorghums have greater biomass than maize and the maximum biomass occurs later in the season. No outliers have been detected at this point. Without a statistical analysis, however, is difficult to make sound statements based solely on data visualization.

**Choose candidate model:** Danalatos et al. (2009) analyzed the data using the beta growth function (Yin et al., 2003a; Eq. [2.5] in Table 1). That model was selected because it captures the decline of biomass toward the end of the growing season (Fig. 4 and supplementary figure for the beta growth function). Also, the parameters have clear meaning and are very suitable to answer the research questions.

**Starting values:** Because for this function the parameters have a straightforward interpretation, starting values can be determined by visual inspection of Fig. 4. In this example, however, we used a prespecified algorithm that chooses the initial starting values automatically (see details in supplementary materials).

**Fit model and convergence:** Model fit was performed in the R package using the ordinary least squares estimation method (nls function). There are three crops, two levels of agronomic input, and four blocks, which results in 24 possible combinations (experimental units). The model was fitted to every experimental unit separately, and apparent convergence was obtained for only 10 experiment units. This indicates that some modifications are needed (see below). Checking model assumptions can be useful for diagnosing the problem (Fig. 5). In this case,

Table 2. Estimates of the beta growth model (Eq. [2.11] in Supplementary Table S2) used to fit the biomass data reported by Danalatos et al. (2009); **P** values < 0.05 indicate a significant effect of input levels (high or low). Note: in Eq. [2.11] the parameters biomass weight at sowing ( $Y_b$ ) and sowing date ( $t_b$ ) were fixed at 0 Mg ha<sup>-1</sup> and Day of the Year (DOY) 141, respectively.

Parameter†	Maize			Fiber sorghum			Sweet sorghum		
	High	Low	<b>P</b>	High	Low	<b>P</b>	High	Low	<b>P</b>
$Y_{\text{max}}$	21.2 (0.99)‡	15.4 (2.27)	<0.00	38.6 (2.24)	31.8 (5.18)	0.02	43.2 (2.83)	33.9 (6.48)	0.01
$t_m$	215.7 (1.30)	217.1 (3.33)	0.50	234.5 (1.61)	235.8 (4.11)	0.61	239.4 (1.53)	240.0 (3.81)	0.79
$t_e$	248.0 (1.79)	248.8 (4.51)	0.76	277.2 (2.03)	279.4 (5.38)	0.50	278.6 (1.99)	279.0 (4.87)	0.89

†  $Y_{\text{max}}$ , maximum biomass (Mg ha<sup>-1</sup>);  $t_m$ , DOY when the crop growth rate is maximized;  $t_e$ , DOY when biomass is maximized.

‡ Standard errors in parentheses.

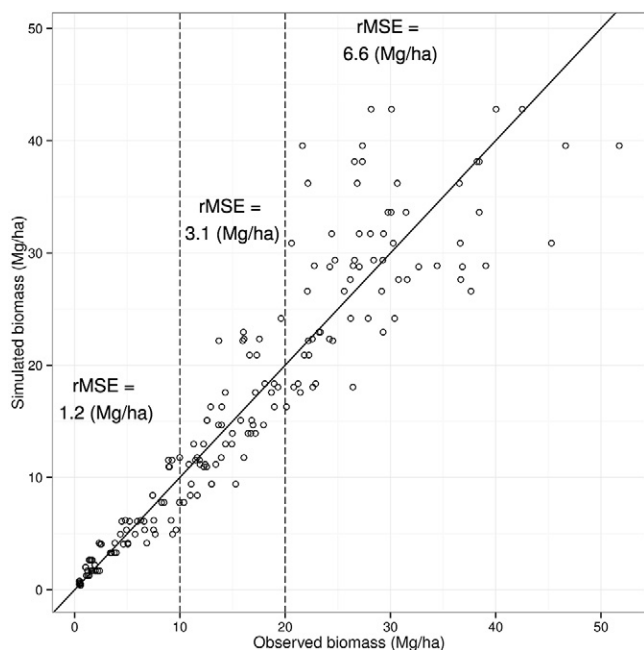


Fig. 6. Observed vs. predicted biomass values. The root mean squared error (RMSE) was used as a measure of the goodness of fit. Given that the variability is increasing along with the biomass weight, three RMSE values were calculated for biomass ranges indicated by the vertical dashed lines (0–10, 10–20, and >20 Mg ha<sup>-1</sup>).

it stands out that there is a concentration of points at low fitted values, which indicates overprediction (i.e., bias) at low values (Fig. 5), suggesting that a different function might work better.

**Revise the mean model:** We selected a modified beta growth function (see Eq. [2.11] in Supplementary Table S2), which was designed to capture more efficiently the initial growth phase at the cost of two extra parameters. Equation [2.11] allows an offset in the  $x$  axis ( $t_b$ ) orientation and an offset in the  $y$  axis orientation ( $Y_b$ ). We did not fit these parameters but rather kept them fixed;  $t_b$  is the planting date at Day of the Year (DOY) 141 and  $Y_b$  is the biomass weight at sowing, which is zero. As a first step, the fitting process was repeated as above,

with starting values determined visually this time as 30 for  $Y_{\max}$ , 240 for  $t_c$ , and 280 for  $t_m$  (Fig. 4). Apparent convergence was obtained for all the experiment units. The final revised mean model was fitted to the entire data set, but at this step the model included the effect of crop type, agronomic input level, and the interaction for each parameter.

**Check model assumptions:** Visual inspection of the standardized residuals (Fig. 5) was used to evaluate the assumptions of appropriate mean function and normally distributed errors with homogeneous variance. Figure 5 indicates that Eq. [2.11] (modified beta growth function) alleviated the overprediction at low values, but this bias did not disappear completely. The major argument for choosing Eq. [2.11] over Eq. [2.5] is that convergence was achieved for all experimental units (24 vs. only 10).

**Model variance homogeneity:** The residual variance was modeled with a power function, and different power parameters were used for the three crops. This function is  $s^2(v) = s^2|v|^{2\theta}$ , where  $v$  is the variance covariate (the fitted values in this case) and  $\theta$  depends on the crop (0.7, 0.86, and 0.89 for maize, fiber sorghum, and sweet sorghum, respectively). More details about the fitting process can be obtained from the supplemental material.

**Determine parameter estimates and standard errors:** Table 2 provides the estimates and the corresponding standard errors of the model. These values are final and account for modeling the residual variance.

**Calculate statistical descriptors:** Given that the biomass had low initial values and high values at the end of the season (Fig. 6), the use of the average RMSE (here 4.1 Mg ha<sup>-1</sup>) is misleading because it overestimates the error at the initial stages (biomass of 0–10 Mg ha<sup>-1</sup>), and it underestimates it at advanced stages (biomass of 30–40 Mg ha<sup>-1</sup>). Therefore, different RMSE values were calculated for different biomass ranges (see Fig. 6). Regarding the relative indices (no units), use of the modeling efficiency (viz. 0.88, scale 0–1) is somewhat better than the RMSE in this case, but it still expresses the average model performance across the sample and therefore is not recommended.

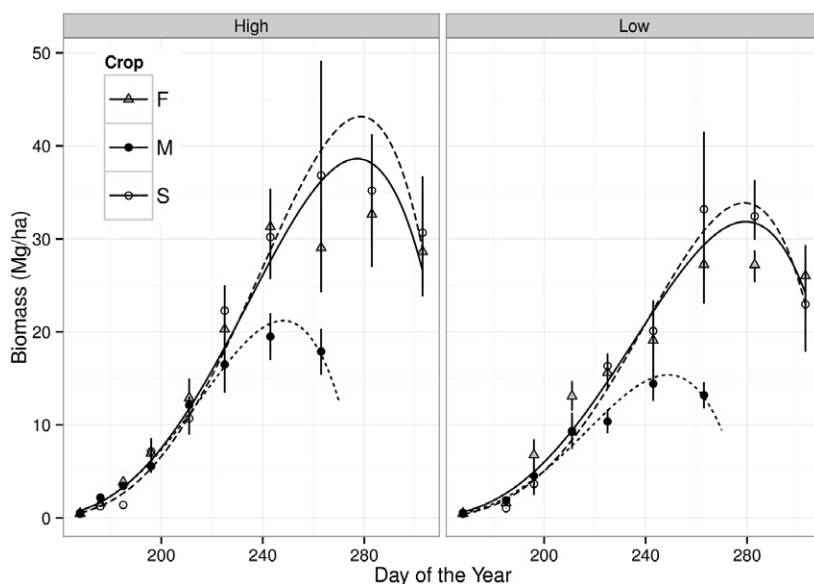


Fig. 7. Observed data and fit for the final model for three crops: maize (M), fiber sorghum (F), and sweet sorghum (S). Vertical bars indicate confidence intervals of observations.

**Interpret results and draw conclusions:** According to the model predictions, the maximum estimated biomass was obtained for sweet sorghum under high inputs, and this crop reached a total of 43 Mg ha<sup>-1</sup> on DOY 279 (Fig. 7; Table 2). At the other extreme, maize reached its maximum biomass under high inputs of 21 Mg ha<sup>-1</sup> on DOY 248. The maximum biomass ( $Y_{\max}$ ) and the time when it was reached ( $t_c$ ) were significantly affected by the crop  $\times$  input interaction (see supplemental materials). In practice, the most meaningful result might be in accurately representing treatment differences and their significance level ( $P$  value) and having a model capable of producing robust predictions within the range of observed values (i.e., interpolation) and, with more caution, outside the range of observed values (i.e., extrapolation).

## SUMMARY

The most critical step that distinguishes nonlinear models from linear models is that the choice of the main function is critical and this can be difficult without appropriate guidance. We have presented an extensive library of nonlinear functions (77 equations with the associated parameter definitions) and typical applications that, we hope, will make the task of choosing candidate models easier. Our review of nonlinear equations is incomplete because there are countless numbers of potential functions (Ratkowsky, 1990) to be used and ad hoc modifications. We have also contributed a suggested work flow (Fig. 1) that should provide the necessary structure to avoid common errors in the use of nonlinear regression models.

## ACKNOWLEDGMENTS

We would like to thank Ken Moore, Philip Dixon, and an anonymous reviewer for their helpful comments and suggestions that improved the manuscript.

## REFERENCES

- Amaducci, S., M. Colauzzi, G. Bellocchi, and G. Venturi. 2008. Modelling post-emergent hemp phenology (*Cannabis sativa* L.): Theory and evaluation. *Eur. J. Agron.* 28:90–102. doi:10.1016/j.eja.2007.05.006
- Andren, O., and K. Paustian. 1987. Barley straw decomposition in the field: A comparison of models. *Ecology* 68:1190–1200. doi:10.2307/1939203
- Archontoulis, S.V., P.C. Struik, X. Yin, L. Bastiaans, J. Vos, and N.G. Danalatos. 2010. Inflorescence characteristics, seed composition, and allometric relationships predicting seed yield in the biomass crop *Cynara cardunculus*. *GCB Bioenergy* 2:113–129. doi:10.1111/j.1757-1707.2010.01045.x
- Archontoulis, S.V., X. Yin, J. Vos, N.G. Danalatos, and P.C. Struik. 2012. Leaf photosynthesis and respiration of three bioenergy crops in relation to temperature and leaf nitrogen: How conservative are biochemical model parameters among crop species? *J. Exp. Bot.* 63:895–911. doi:10.1093/jxb/err321
- Arrhenius, S. 1889. Ueber die Reaktionsgeschwindigkeit bei der Inversion von Rohrzucker durch Säuren. *Z. Phys. Chem.* 4:226–248.
- Atkin, O.K., D. Bruhn, V.M. Hurry, and M.G. Tjoelker. 2005. The hot and the cold: Unraveling the variable response of plant respiration to temperature. *Funct. Plant Biol.* 32:87–105. doi:10.1071/FP03176
- Bates, D.M., and D.G. Watts. 2007. Nonlinear regression analysis and its applications. Wiley Ser. Probab. Stat. John Wiley and Sons, New York.
- Birch, C.P.D. 1999. A new generalized logistic sigmoid growth equation compared with the Richards growth equation. *Ann. Bot.* 83:713–723. doi:10.1006/anbo.1999.0877
- Blackman, F.F. 1905. Optima and limiting factors. *Ann. Bot.* 19:281–295.
- Bolker, B.M. 2008. Ecological models and data in R. Princeton Univ. Press, Princeton, NJ.
- Bril, J., H.G. Van Faassen, and H. Klein Gunnewiek. 1994. Modeling N<sub>2</sub>O emission from grazed grassland. Rep. 24. DLO Res. Inst. for Agrobiol. and Soil Fert., Haren, the Netherlands.
- Burnham, K.P., and D.R. Anderson. 2002. Model selection and multimodel inference: A practical information–theoretical approach. 2nd ed. Springer-Verlag, New York.
- Carroll, R.J., and D. Ruppert. 1988. Transformations and weighting in regression. Chapman and Hall, New York.
- Collatz, G.J., M. Ribas-Carbo, and J.A. Berry. 1992. Coupled photosynthesis–stomatal conductance model for leaves of C<sub>4</sub> plants. *Aust. J. Plant Physiol.* 19:519–538. doi:10.1071/PP9920519
- Danalatos, N.G., S.V. Archontoulis, and K. Tsiboukas. 2009. Comparative analysis of sorghum versus corn growing under optimum and under water/nitrogen limited conditions in central Greece. In: From research to industry and markets: Proceedings of the 17th European Biomass Conference, Hamburg, Germany. 29 June–3 July 2009. ETA–Renewable Energies, Florence, Italy. p. 538–544.
- Davidson, E.A., I.A. Janssens, and Y. Luo. 2006. On the variability of respiration in terrestrial ecosystems: Moving beyond Q<sub>10</sub>. *Global Change Biol.* 12:154–164. doi:10.1111/j.1365-2486.2005.01065.x
- Farquhar, G.D., S. von Caemmerer, and J.A. Berry. 1980. A biochemical model of photosynthetic CO<sub>2</sub> assimilation in leaves of C<sub>3</sub> species. *Planta* 149:78–90. doi:10.1007/BF00386231
- Gan, Y., E.H. Stobbe, and C. Njue. 1996. Evaluation of selected nonlinear regression models in quantifying seedling emergence rate of spring wheat. *Crop Sci.* 36:165–168. doi:10.2135/cropsci1996.0011183X003600010029x
- Gbur, E.E., W.W. Stroup, K.S. McCarter, S. Durham, L.J. Young, M. Christman, M. West, and M. Kramer. 2012. Analysis of generalized linear mixed models in the agricultural and natural resources sciences. ASA, CSSA, and SSSA, Madison, WI.
- Gillis, J.D., and G.W. Price. 2011. Comparison of a novel model to three conventional models describing carbon mineralization from soil amended with organic residues. *Geoderma* 160:304–310. doi:10.1016/j.geoderma.2010.09.025
- Gompertz, B. 1825. On the nature of the function expressive of the law of human mortality, and on a new mode of determining the value of life contingencies. *Philos. Trans. R. Soc.* 115:513–585. doi:10.1098/rstl.1825.0026
- Goudriaan, J. 1979. A family of saturation type curves, especially in relation to photosynthesis. *Ann. Bot.* 43:783–785.
- Goudriaan, J., and H.H. van Laar. 1994. Modelling potential crop growth processes. Kluwer Acad. Publ., Dordrecht, the Netherlands.
- Hammer, G.L., Z. Dong, G. McLean, A. Doherty, C. Messina, J. Schussler, et al. 2009. Can changes in canopy and/or root system architecture explain historical maize yield trends in the U.S. Corn Belt? *Crop Sci.* 49:299–312. doi:10.2135/cropsci2008.03.0152
- Hoben, J.P., R.J. Gehl, N. Millar, P.R. Grace, and G.P. Robertson. 2011. Nonlinear nitrous oxide (N<sub>2</sub>O) response to nitrogen fertilizer in on-farm corn crops of the US Midwest. *Global Change Biol.* 17:1140–1152. doi:10.1111/j.1365-2486.2010.02349.x
- Hoffmann, J.P. 2005. Linear regression analysis: Assumptions and applications. SPSS version. Dep. of Sociology, Brigham Young Univ., Provo, UT. [https://sociology.byu.edu/Hoffmann/SiteAssets/Hoffmann%20\\_%20Linear%20Regression%20Analysis.pdf](https://sociology.byu.edu/Hoffmann/SiteAssets/Hoffmann%20_%20Linear%20Regression%20Analysis.pdf) (accessed 25 Apr. 2013).
- Horton, R.E. 1940. An approach towards a physical interpretation of infiltration capacity. *Soil Sci. Soc. Am. Proc.* 5:227–237.
- Huet, S., A. Bouvier, M.-A. Poursat, and E. Jolivet. 2003. Statistical tools for nonlinear regression: A practical guide with S-PLUS and R examples. 2nd ed. Springer Ser. Stat. Springer-Verlag, New York.
- Jassby, A.D., and T. Platt. 1976. Mathematical formulation of the relationship between photosynthesis and light for phytoplankton. *Limnol. Oceanogr.* 21:540–547. doi:10.4319/lo.1976.21.4.0540
- Johnson, I.R., J.H.M. Thornley, J.M. Frantz, and B. Bugbee. 2010. A model of canopy photosynthesis incorporating protein distribution through the canopy and its acclimation to light, temperature and CO<sub>2</sub>. *Ann. Bot.* 106:735–749. doi:10.1093/aob/mcq183
- Kätterer, T., M. Reichstein, O. Andrén, and A. Lomander. 1998. Temperature dependence of organic matter decomposition: A critical review using literature data analyzed with different models. *Biol. Fertil. Soils* 24:258–262. doi:10.1007/s003740050430
- Keating, B.A., P.S. Carberry, G.L. Hammer, M.E. Probert, M.J. Robertson, D. Holzworth, et al. 2003. An overview of APSIM, a model designed for farming systems simulation. *Eur. J. Agron.* 18:267–288. doi:10.1016/S1161-0301(02)00108-9



- Kirschbaum, M.U.F. 1995. The temperature dependence of soil organic matter decomposition, and the effect of global warming on soil C storage. *Soil Biol. Biochem.* 27:753–760. doi:10.1016/0038-0717(94)00242-S
- Leduc, D., and L. Goetz. 2009. A height–diameter curve for longleaf pine plantations in the Gulf Coastal Plain. *South. J. Appl. For.* 33:164–170.
- Lewis, F., A. Butler, and L. Gilbert. 2011. A unified approach to model selection using the likelihood ratio test. *Methods Ecol. Evol.* 2:155–162. doi:10.1111/j.2041-210X.2010.00063.x
- Lindsey, J.K. 2001. *Nonlinear models for medical statistics*. 2nd ed. Oxford Stat. Sci. Ser. 24. Oxford Univ. Press, Oxford, UK.
- Lloyd, J., and J.A. Taylor. 1994. On the temperature dependence of soil respiration. *Funct. Ecol.* 8:315–323. doi:10.2307/2389824
- Ma, L., and M.J. Shaffer. 2001. A review of carbon and nitrogen processes in nine U.S. soil nitrogen dynamic models. In: M.J. Shaffer et al., editors, *Modeling carbon and nitrogen dynamics for soil management*. Lewis Publ., Boca Raton, FL. p. 55–103.
- McGeachan, M.B., and L. Wu. 2001. A review of carbon and nitrogen processes in European soil nitrogen dynamic models. In: M.J. Shaffer et al., editors, *Modeling carbon and nitrogen dynamics for soil management*. Lewis Publ., Boca Raton, FL. p. 103–171.
- Medlyn, B.E., E. Dreyer, D. Ellsworth, M. Forstreuter, P.C. Harley, M.U.F. Kirschbaum, et al. 2002. Temperature response of parameters of a biochemically based model of photosynthesis: II. A review of experimental data. *Plant Cell Environ.* 25:1167–1179. doi:10.1046/j.1365-3040.2002.00891.x
- Miguez, F.E., M.B. Villamil, S.P. Long, and G.A. Bollero. 2008. Meta-analysis of the effects of management factors on *Miscanthus × giganteus* growth and biomass production. *Agric. For. Meteorol.* 148:1280–1292. doi:10.1016/j.agrformet.2008.03.010
- Miguez, F.E., X. Zhu, S. Humphries, G.A. Bollero, and S.P. Long. 2009. A semi-mechanistic model predicting growth and production of the bioenergy crop *Miscanthus × giganteus*: Description, parameterization and validation. *GCB Bioenergy* 1:282–296. doi:10.1111/j.1757-1707.2009.01019.x
- Monsi, M., and T. Saeki. 2005. On the factor light in plant communities and its importance for matter production. *Ann. Bot.* 95:549–567. doi:10.1093/aob/mci052
- Myung, I.J. 2003. Tutorial on maximum likelihood estimation. *J. Math. Psychol.* 47:90–100. doi:10.1016/S0022-2496(02)00028-7
- Osório, M.L., E. Breia, A. Rodrigues, J. Osório, X. Le Roux, F.A. Daudet, et al. 2006. Limitations to carbon assimilation by mild drought in nectarine trees growing under field conditions. *Environ. Exp. Bot.* 55:235–247. doi:10.1016/j.envexpbot.2004.11.003
- Pinheiro, J.C., and D.M. Bates. 2000. *Mixed-effects models in S and S-PLUS*. Stat. Comput. Ser. Springer-Verlag, New York.
- Portner, H., H. Bugmann, and A. Wolf. 2010. Temperature response functions introduce high uncertainty in modelled carbon stocks in cold temperature regimes. *Biogeosciences* 7:3669–3684. doi:10.5194/bg-7-3669-2010
- Rafigue, R. 2011. *Measurements and modeling of nitrous oxide emissions from Irish grasslands*. Ph.D. diss. Natl. Univ. of Ireland, Cork.
- Ratkowsky, D.A. 1990. *Handbook of nonlinear regression models*. Marcel Dekker, New York.
- Ratkowsky, D.A., J. Olley, T.A. McMeekin, and A. Ball. 1982. Relationship between temperature and growth rate of bacterial cultures. *J. Bacteriol.* 149:1–5.
- Richards, F.J. 1959. A flexible growth function for empirical use. *J. Exp. Bot.* 10:290–300. doi:10.1093/jxb/10.2.290
- Ritz, C., and J.C. Streibig. 2005. *Bioassay analysis using R*. J. Stat. Softw. 12:1–21.
- Ritz, C., and J.C. Streibig. 2008. *Nonlinear regression with R*. Springer, New York.
- Ruppert, D., N. Cressie, and R.J. Carroll. 1989. A transformation/weighting model for estimating Michaelis–Menten parameters. *Biometrics* 45:637–656. doi:10.2307/2531506
- Shibu, M.E., P.A. Leffelaar, H. van Keulen, and P.K. Aggarwal. 2006. Quantitative description of soil organic matter dynamics: A review of approaches with reference to rice-based cropping systems. *Geoderma* 137:1–18. doi:10.1016/j.geoderma.2006.08.008
- Sinclair, T.R., and T. Horie. 1989. Leaf nitrogen, photosynthesis, and crop radiation use efficiency: A review. *Crop Sci.* 29:90–98. doi:10.2135/cropsci1989.0011183X002900010023x
- Singh, P. 2006. *Modeling crop production systems: Principles and application*. Sci. Publ., Enfield, NH.
- Tjoelker, M.G., L. Oleksyn, and P. Reich. 2001. Modelling respiration of vegetation: Evidence for a general temperature-dependent  $Q_{10}$ . *Global Change Biol.* 7:223–230. doi:10.1046/j.1365-2486.2001.00397.x
- Tsoularis, A. 2001. Analysis of logistic growth models. *Res. Lett. Inf. Math. Sci.* 2:23–46.
- van't Hoff, J.H. 1898. *Lectures on theoretical and physical chemistry*. Part 1. Chemical dynamics. Edward Arnold, London.
- Vega, C.R.C., and V.O. Sadras. 2003. Size dependent growth and development of inequality in maize, sunflower and soybean. *Ann. Bot.* 91:795–805. doi:10.1093/aob/mcg081
- Vega, C.R.C., V.O. Sadras, F.H. Andrade, and S.A. Uhart. 2000. Reproductive allometry in soybean, maize and sunflower. *Ann. Bot.* 85:461–468. doi:10.1006/anbo.1999.1084
- Verhulst, P.F. 1838. A note on population growth. *Corresp. Math. Phys.* 10:113–121.
- Vico, G., and A. Porporato. 2008. Modelling  $C_3$  and  $C_4$  photosynthesis under water-stressed conditions. *Plant Soil* 313:187–203. doi:10.1007/s11040-008-9691-4
- von Caemmerer, S., and G.D. Farquhar. 1981. Some relationships between the biochemistry of photosynthesis and the gas exchange of leaves. *Planta* 153:376–387. doi:10.1007/BF00384257
- Wallach, D. 2006. Evaluating crop models. In: D. Wallach et al., editors, *Working with dynamic crop models: Evaluations, analysis, parameterization, and applications*. Elsevier, Amsterdam. p. 11–53.
- Wang, E., and T. Engel. 1998. Simulation of phenological development of wheat crops. *Agric. Syst.* 58:1–24. doi:10.1016/S0308-521X(98)00028-6
- Weibull, W. 1951. A statistical distribution function of wide applicability. *J. Appl. Math.* 18:293–297.
- Yang, H.S., A. Dobermann, J.L. Lindquist, D.T. Walters, T.J. Arkebauer, and K.G. Cassman. 2004. Hybrid-maize: A maize simulation model that combines two crop modeling approaches. *Field Crops Res.* 87:131–154. doi:10.1016/j.fcr.2003.10.003
- Yin, X., J. Goudriaan, E.A. Lantinga, J. Vos, and J.H.J. Spiertz. 2003a. A flexible sigmoid function of determinate growth. *Ann. Bot.* 91:361–371; 753 (with erratum in *Annals of Botany* 91:753, 2003).
- Yin, X., M.J. Kroff, G. McLean, and R.M. Visperas. 1995. A nonlinear model for crop development as a function of temperature. *Agric. For. Meteorol.* 77:1–16. doi:10.1016/0168-1923(95)02236-Q
- Yin, X., E.A. Lantinga, A.H.C.M. Schapendonk, and X. Zhong. 2003b. Some quantitative relationships between leaf area index and canopy nitrogen content and distribution. *Ann. Bot.* 91:893–903. doi:10.1093/aob/mcg096
- Yin, X., A.H.C.M. Schapendonk, M.J. Kropff, M. van Oijen, and P.S. Bindra-ban. 2000. A generic equation for nitrogen-limited leaf area index and its application in crop growth models for predicting leaf senescence. *Ann. Bot.* 85:579–585. doi:10.1006/anbo.1999.1104
- Yin, X., and P.C. Struik. 2009.  $C_3$  and  $C_4$  photosynthesis models: An overview from the perspective of crop modelling. *NJAS–Wageningen J. Life Sci.* 57:27–38. doi:10.1016/j.njas.2009.07.001
- Zeide, B. 1993. Analysis of growth equations. *For. Sci.* 39:594–616.
- Zucchini, W. 2000. An introduction to model selection. *J. Math. Psychol.* 44:41–61. doi:10.1006/jmps.1999.1276
- Zwietering, M.H., I. Jongenburger, F.M. Rombouts, and K. van't Riet. 1990. Modeling of the bacterial growth curve. *Appl. Environ. Microbiol.* 56:1875–1881.