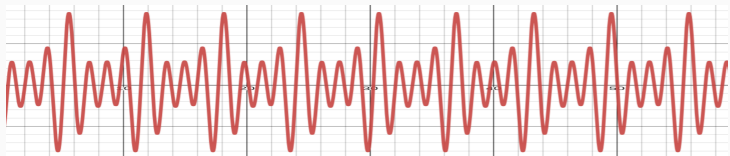


Model Validation



Prof. Fabrício Olivetti de França

Federal University of ABC

05 February, 2024



Model Validation

As already stressed throughout this course, there are three main approaches for nonlinear regression:

- Using an overparameterized generic model (opaque model).
- Manually crafting the nonlinear model.
- Using Symbolic Regression to find a nonlinear model with as few parameters as possible.

While crafting the model using first principles, you may have some properties that you want to enforce into your model, either because of some requirements or from a prior knowledge about the behavior of the system.

In this situation, the practitioner can enforce those using their own expertise.

For example, due to EU regulations¹, the practitioner will create a model that will allow them to debug how the output is generated in a clear manner. Also, they may want to ensure fairness in the predictions.

¹(<https://www.europarl.europa.eu/news/en/press-room/20240308IPR19015/artificial-intelligence-act-meps-adopt-landmark-law>)[<https://www.europarl.europa.eu/news/en/press-room/20240308IPR19015/artificial-intelligence-act-meps-adopt-landmark-law>]

This is usually a problem for opaque models that are often hard to debug and not flexible enough to enforce some properties of interest.

In the current literature, there are some techniques that can extract information from opaque models to have a better understanding. But this may not be enough in practice.

With the *vanilla* symbolic regression, you have the possibility of finding a model that attends to all your requirements. To increase the probability of finding the correct model, you need at least one of these:

- Noiseless data.
- Representative data.
- Luck 🍀
- A well calibrated SR algorithm.

With the *vanilla* symbolic regression, you have the possibility of finding a model that attends to all your requirements. To increase the probability of finding the correct model, you need at least one of these:

- Noiseless data.
- Representative data.
- Luck 🍀
- A well calibrated SR algorithm.

We can only afford the last one!

Another important motivation for model validation is that, depending on the hyper-parameters, the SR algorithm can favor large and overparameterized models that will have a high goodness-of-fit without the remaining desiderata.

Some example of objectives beyond the goodness-of-fit² are:

- The ability to understand and explain model behavior
- Scientific plausibility of the model
- Whether the model is generalizable and capable of extrapolation
- Boundedness and safe operation under all circumstances
- Efficiency of calculating predictions or computational effort required for training the model

²Gabriel Kronberger, Bogdan Burlacu, Michael Kommenda, Stephan M. Winkler, and Michael Affenzeller. Symbolic Regression. tbr.

Besides those, we may also want a model that:

- Ensures a fair inference to different classes of the sample.
- Behaves according to pre-established norms.

In the beginning of the course, it was clear that a linear model is easy to understand:

- With every unitary change in x we observe a change proportional to β in the outcome.
- Even if we have a linear model with non-linear features, they can have physical meaning. E.g., $v = s/t$, the inverse interaction of displacement and time gives us the average velocity.

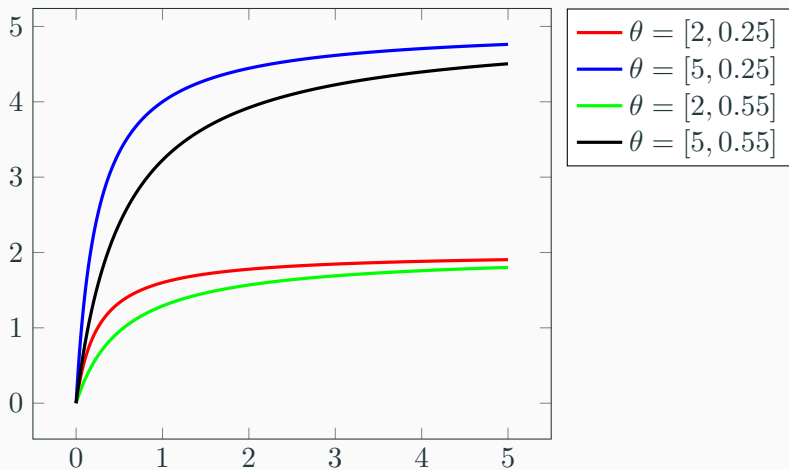
When we have a nonlinear regression model, these interpretations are not as straightforward:

$$f(x; \theta) = \frac{\theta_1 x}{\theta_2 + x},$$

The association between the input variable and the outcome is not easily understood.

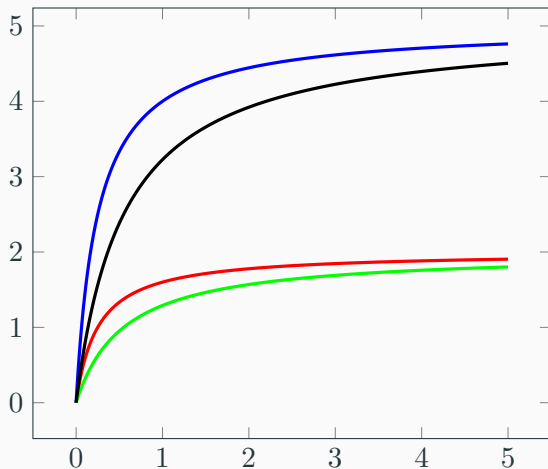
Ability to understand and explain model behavior

We can try to understand the behavior with a plot for different values of θ :



Ability to understand and explain model behavior

::::column ::column



::: ::column - This model has a saturation value close to θ_1 - The higher the value of θ_1 , the longer the model takes to reach the saturation. When $\theta_1 = 0$

Having the context of the model can help gain additional insights. This particular model can represent the **Michaelis–Menten kinetics** that describes the reaction rate ($f(x; \theta)$) to the concentration of a substrate (x).

Knowing the physical meanings of θ will give us insight when fitting this model for different enzymes.

We can see that, once we contextualize the model and add expert knowledge, we can gain insights from nonlinear models as well, as long as their parameters are meaningful in our context (thus, minimize the number of parameters is desired).

In short, inspecting the model for the ability of understanding and explaining can be done by:

- Contextualizing the model
- Applying expert knowledge
- Plotting the behavior of the function with different parameter values

Additional tools will be given in later lectures when we talk about explainability.

Related to the previous desiderata, scientific plausibility refers to whether the model:

- Behaves similarly to the observed phenomena.
- Is correct w.r.t. a dimensional analysis (or whether all meta-features are dimensionless)
- Possesses a physical meaning
- Does not misbehave

This can be inspected through visual plots and expert knowledge.

Whether the model is generalizable and capable of extrapolation

Boundedness and safe operation under all circumstances

Efficiency of calculating predictions or computational effort required for training the model

Ensures a fair inference to different classes of the sample.

Behaves according to pre-established norms.

- Model Selection



Acknowledgments