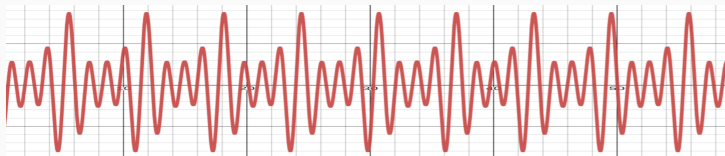


# Model Validation



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 regulations<sup>1</sup>, 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.

---

<sup>1</sup>(<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-fit<sup>2</sup> 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

---

<sup>2</sup>Gabriel 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  $\beta$  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.

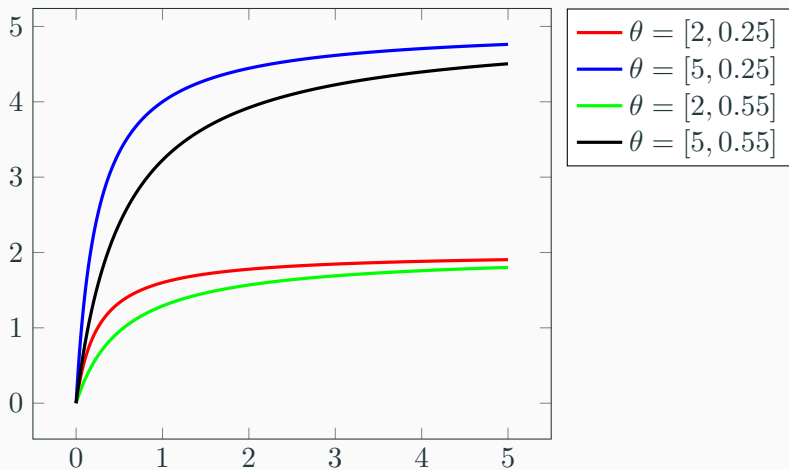
When we have a nonlinear regression model, these interpretations are not as straightforward:

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

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

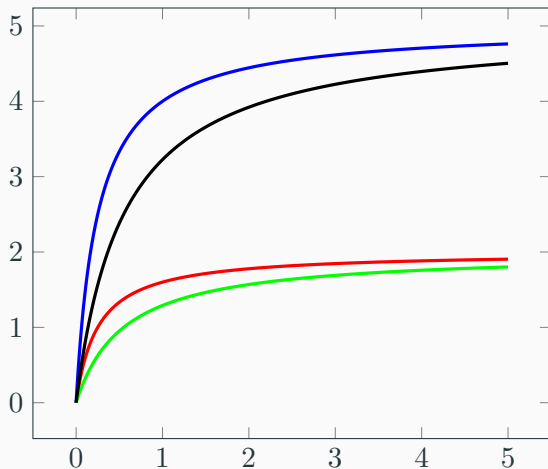
## Ability to understand and explain model behavior

We can try to understand the behavior with a plot for different values of  $\theta$ :



## Ability to understand and explain model behavior

::::column ::column



::::column - This model has a saturation value close to  $\theta_1$  - The higher the value of  $\theta_1$ , 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  $\theta$  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

The SR model is fitted on a limited data set that not necessarily captures a

## 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.**

The plots regarding the predictions and residuals can be insightful and provide a tool for model inspection. From these plots we can understand whether the model meets our expectations and whether there is any unexpected behavior.



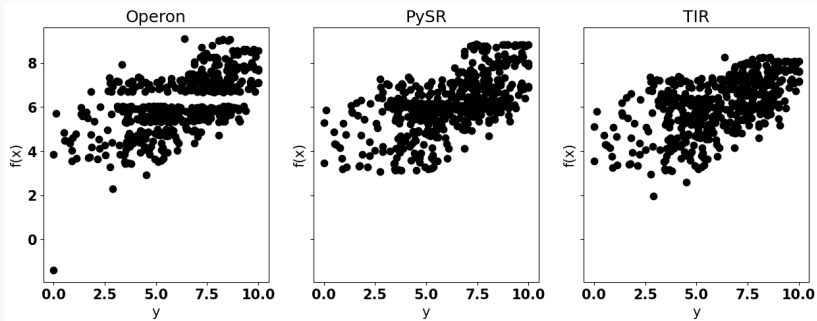
One example of a plot is the predicted values against the dependent variable as observed in the data. To illustrate this and the next plots, we will fit our simulated grade dataset with PyOperon, PySR and TIR:

```
1 regs = [SymbolicRegressor(),  
2         TIRRegressor(100, 100, 0.3, 0.7, (-3, 3), transfunctions='Id',  
3         PySRRegressor(binary_operators=["+", "*"], unary_operators=[])  
4         ]  
5 for i in range(3):  
6     regs[i].fit(x.reshape(-1,1),y)
```

We can check the noise variance with:

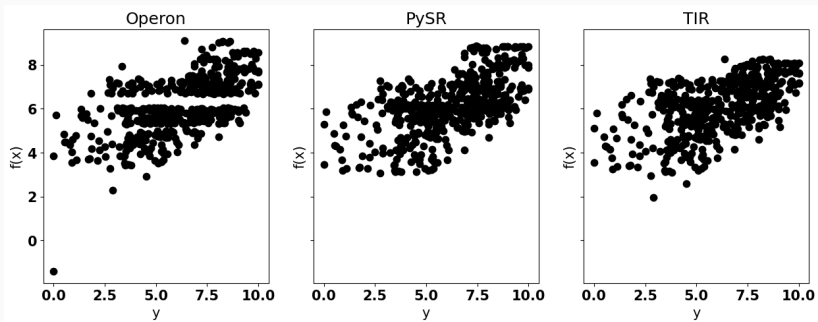
```
1 _,axs = plt.subplots(1,3, figsize=(15,5), sharey=True)
2 name = ['Operon', 'PySR', 'TIR']
3 for i in range(3):
4     axs[i].plot(y, regs[i].predict(x.reshape(-1,1)), '.', color='black',
5     axs[i].set_xlabel('y')
6     axs[i].set_ylabel('f(x)')
7     axs[i].set_title(name[i])
```

## Noise variance plot



A perfect model would have all the points in the 45 degrees diagonal.

## Noise variance plot



We can see from these plots that none of the models returns a satisfactory result. Also, we can see that all of them have a bias in mispredicting grades below 5 (usually for a higher grade).

Another important plot is the quantile-quantile plot (Q-Q plot) that plots the assumed error distribution of the data matches the distribution of the residuals of the model.

To make the Q-Q plot, we calculate the residuals of our model, sort them in increasing order, and plot each point against the inverse of the cumulative density function of the assumed distribution.

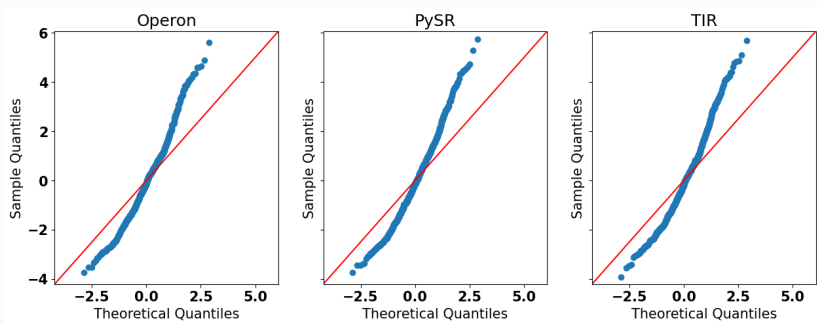
(qqplot assumes normal distribution as the default)

---

```
1 import statsmodels.api as sm
2
3 _,axs = plt.subplots(1,3, figsize=(15,5), sharey=True)
4
5 for i in range(3):
6     sm.qqplot(regs[i].predict(x.reshape(-1,1))[:,0]-y, line = '45', ax=axs[i])
7     axs[i].set_title(name[i])
```

---

## Q-Q plot



We can see from these plots that none of the models matches the expected distribution for the residuals.

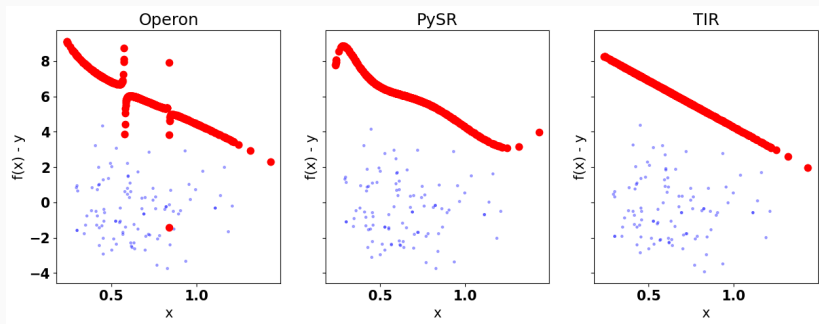
## Residuals plot

Another interesting plot is the residuals plots in which we plot a choice of  $x_i$  against  $f(x)$  and the residuals:

```
1 import statsmodels.api as sm
2
3 _,axs = plt.subplots(1,3, figsize=(15,5), sharey=True)
4
5 for i in range(3):
6     axs[i].plot(x, regs[i].predict(x.reshape(-1,1))[:,0] - y, '.', color=
7     axs[i].plot(x, regs[i].predict(x.reshape(-1,1)), '.', color='red', ma
8     axs[i].set_xlabel('x')
9     axs[i].set_ylabel('f(x) - y')
10    axs[i].set_title(name[i])
```

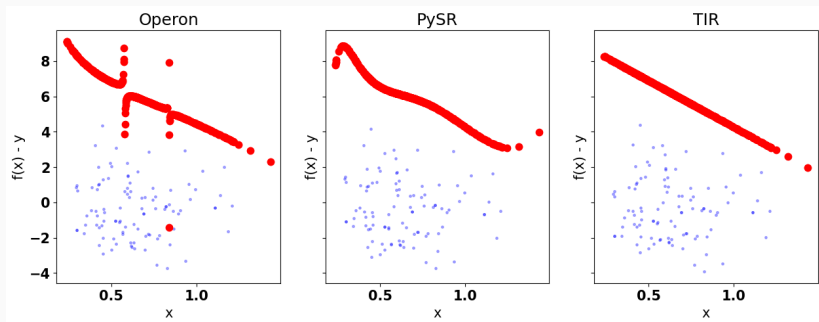


# Residuals plot



These plots show that all of these models have an error ranging from  $-4$  to  $4$  but mostly concentrated on negative residuals. This means it tends to underestimate the true value.

# Residuals plot



Also, we can see that Operon created a model with some discontinuities (possibly because of division) and TIR chose a linear model.

- Model Selection



# Acknowledgments