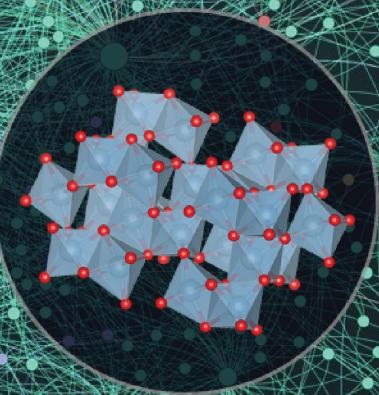
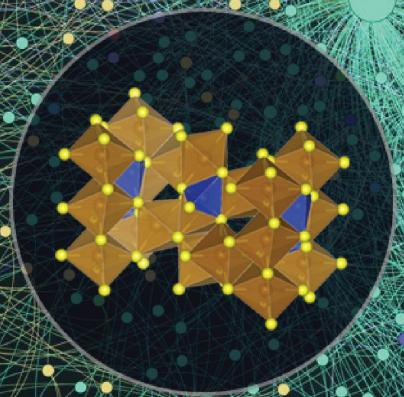
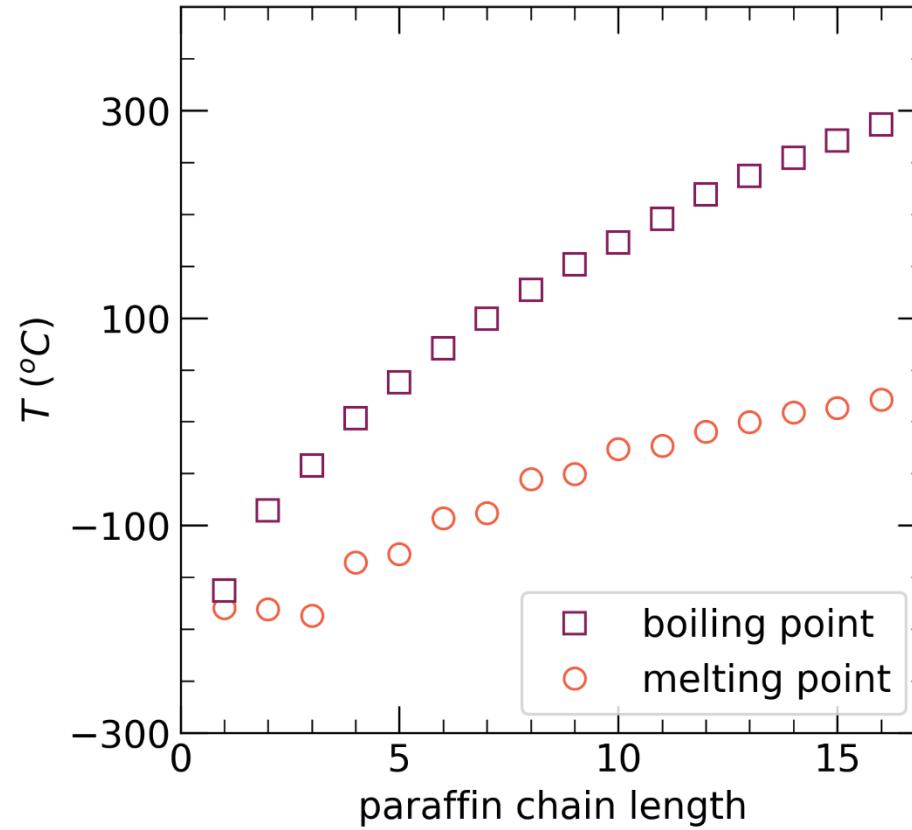
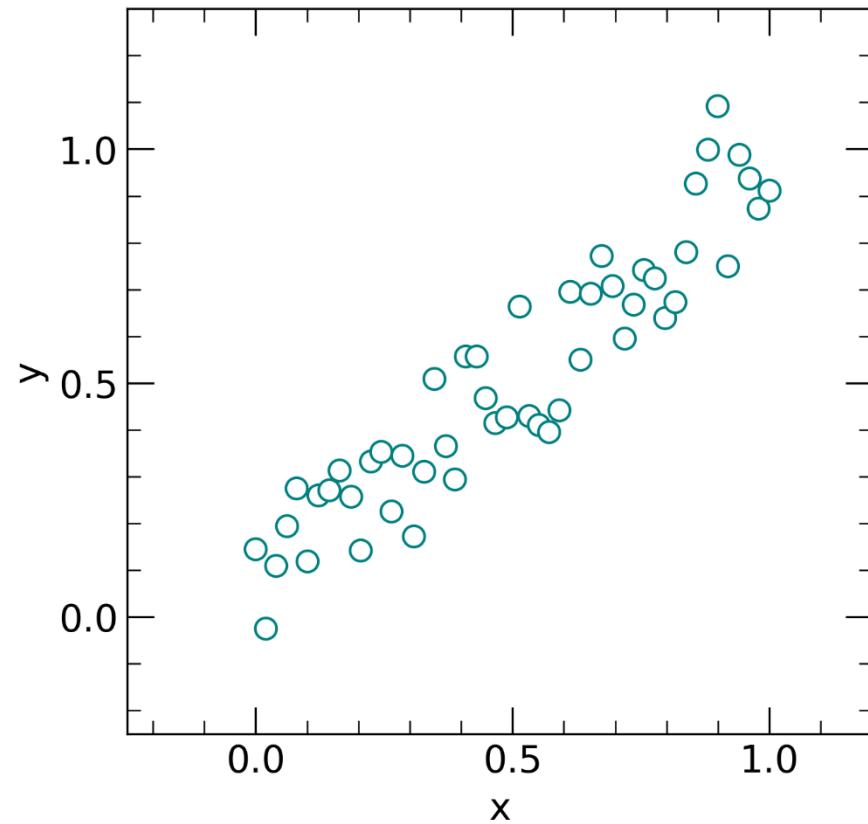


Linear and Nonlinear models



Both linear and non-linear models can fit curvature in data



Linear models must follow a strict format

*dependent variable = constant + parameter * IV + parameter * IV + ...*

Linear models must follow a strict format

*dependent variable = constant + parameter * IV + parameter * IV + ...*

$$y = mx + b$$

$$x_{final} = x_{initial} + v_0 t + 1/2 a t^2$$



Linear models must follow a strict format

*dependent variable = constant + parameter * IV + parameter * IV + ...*

$$y = mx + b$$

$$x_{final} = x_{initial} + v_0 t + 1/2 a t^2$$

Linear models consist of sums of parameters multiplied by independent variables (IV)

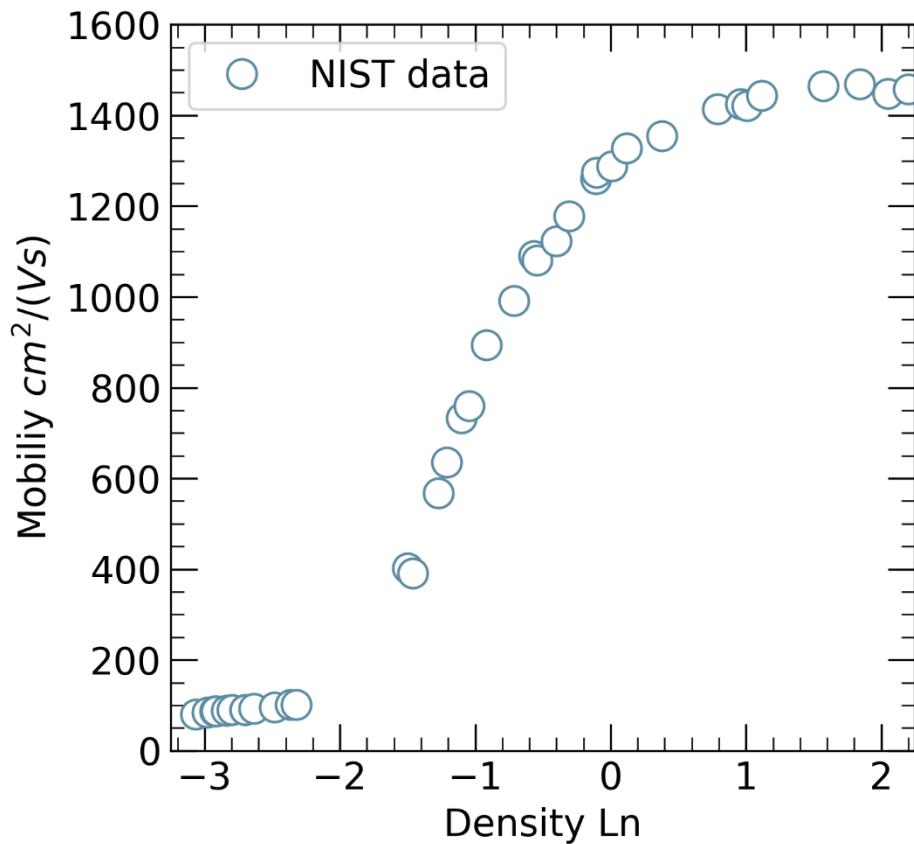
Nonlinear models are anything that don't follow this form



Linear models must follow a strict format

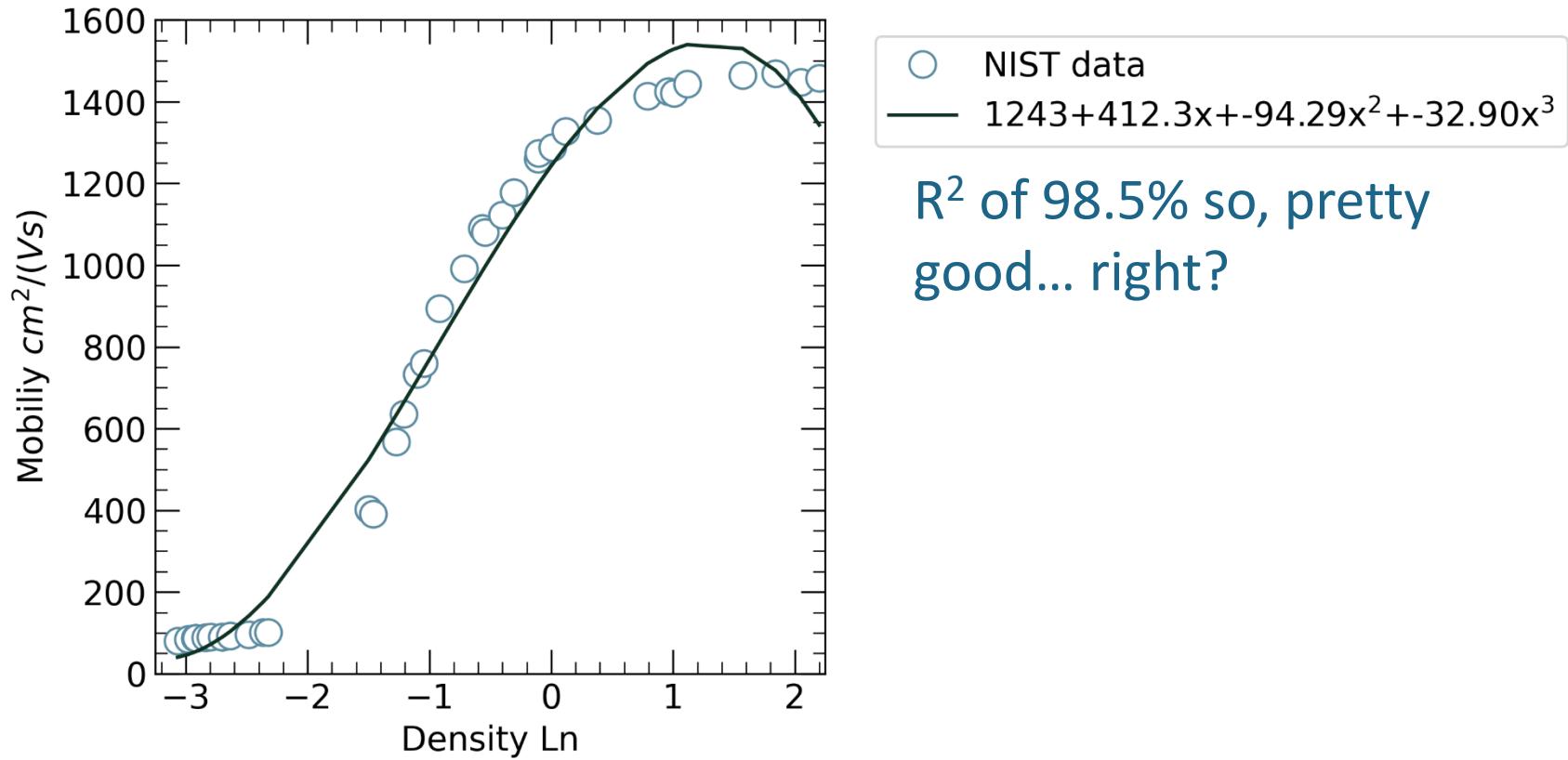
Benefits of linear models	Benefits of nonlinear models
Easy to implement	Harder to implement
Easy to interpret	Harder to interpret
Easily obtained statistics to assess model	Can fit more complicated trends in data (but this can be hard to accomplish)
	No R^2 , no p-values possible for the parameter estimates

Let's look at an example for materials data



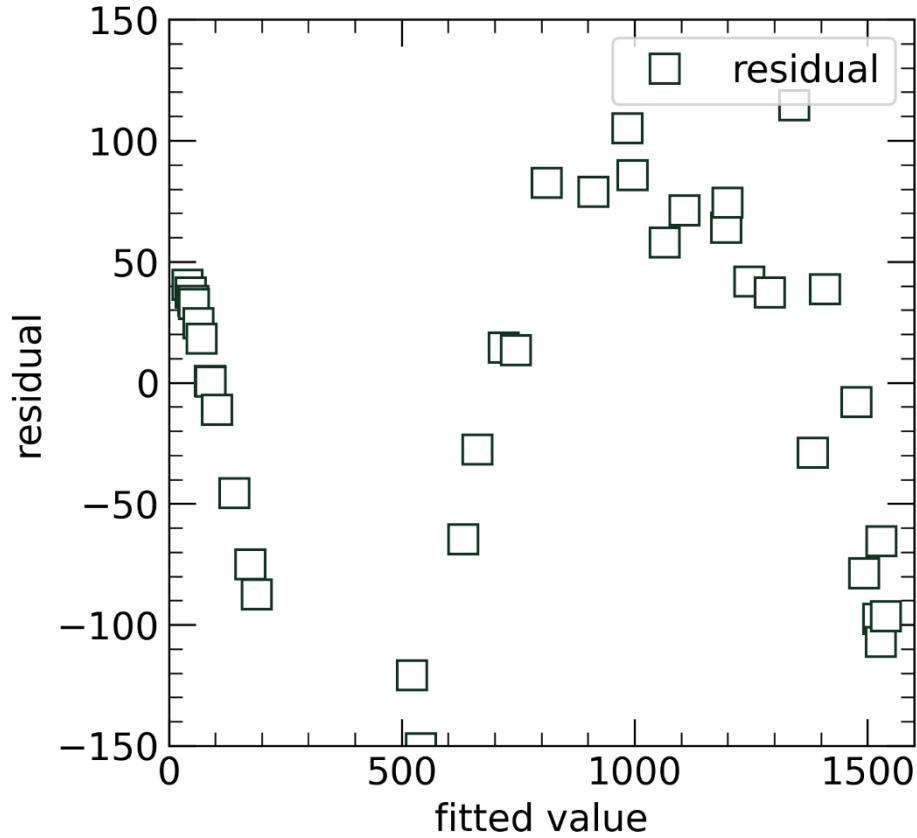
Data taken from NIST <https://www.itl.nist.gov/div898/strd/nls/data/thurber.shtml>

A third order polynomial is an OK fit...



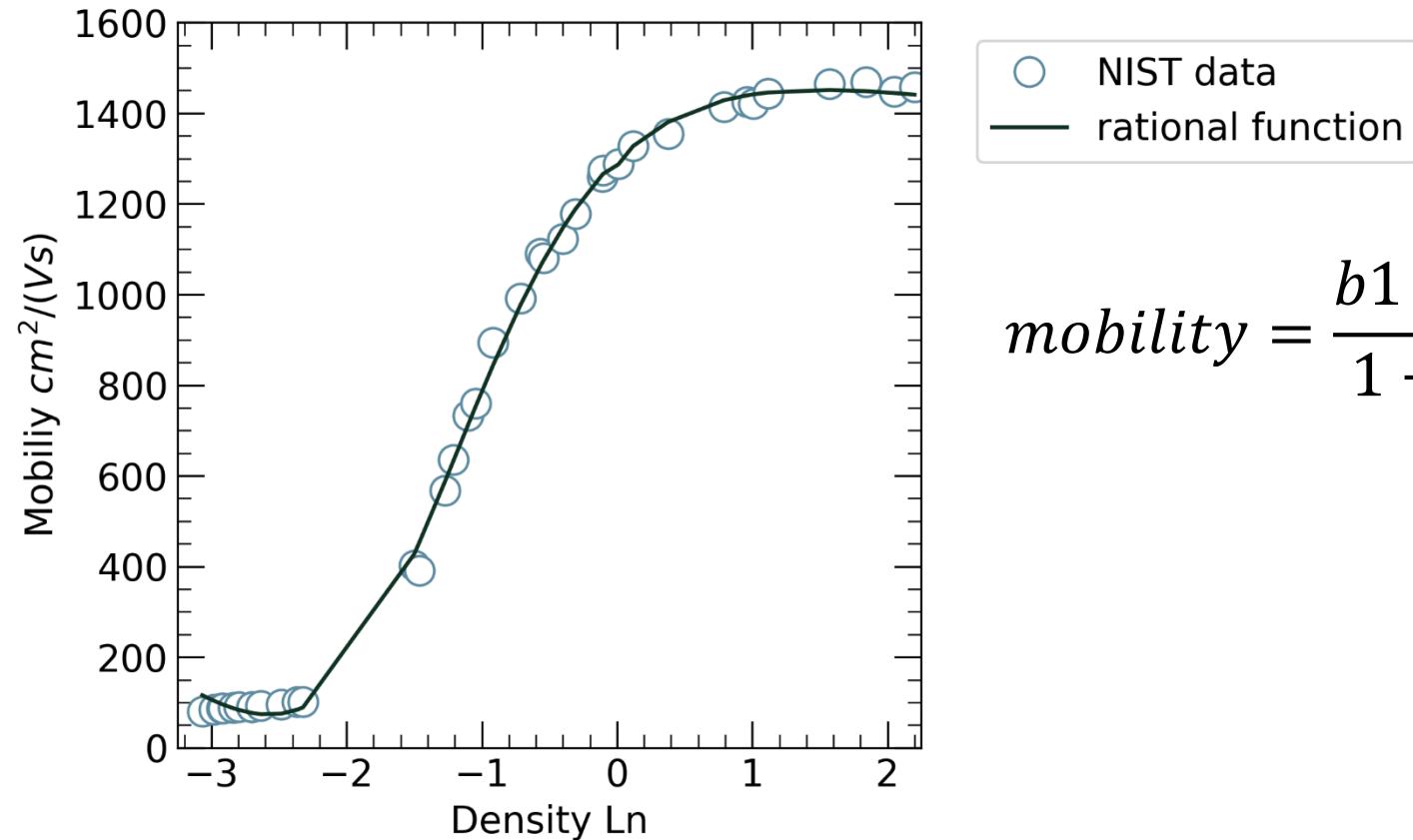
Clearly it overestimates and underestimates portions of the data in systematic (biased) ways

If we examine the residual, we see the bias



A model that is systematically incorrect is biased (R^2 doesn't tell us the whole story!)

Non-linear fit is also possible... but it's trickier!

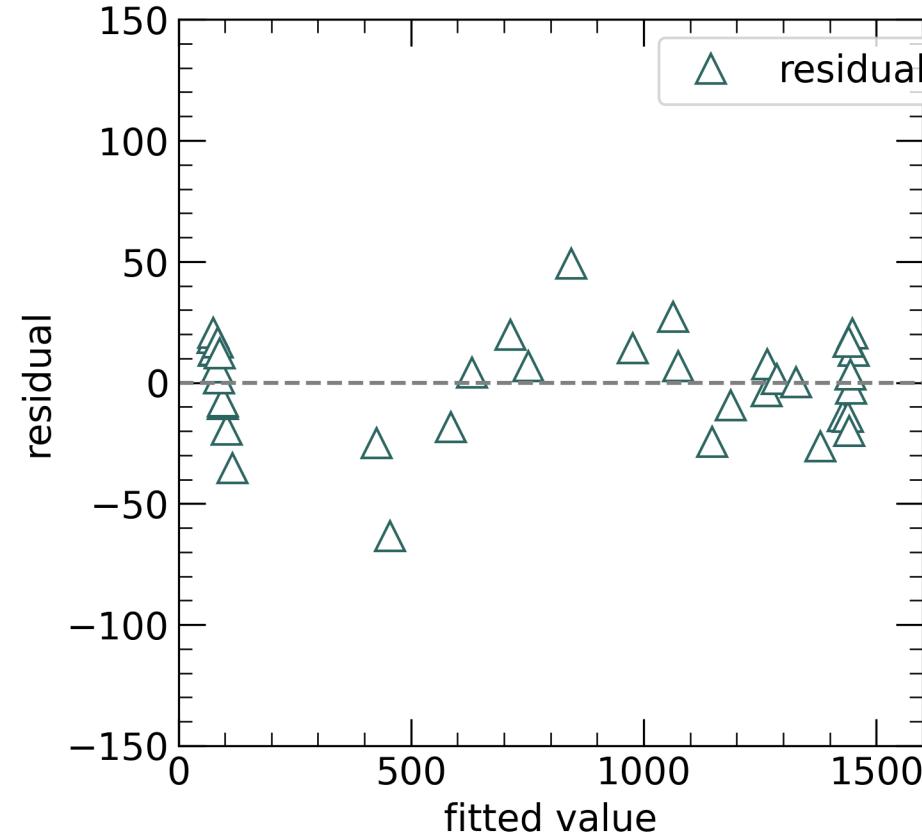


$$mobility = \frac{b_1 + b_2x + b_3x^2 + b_4x^3}{1 + b_5x + b_6x^2 + b_7x^3}$$

A rational function is two polynomials divided by one another
Guessing the starting values for the fit is tricky!

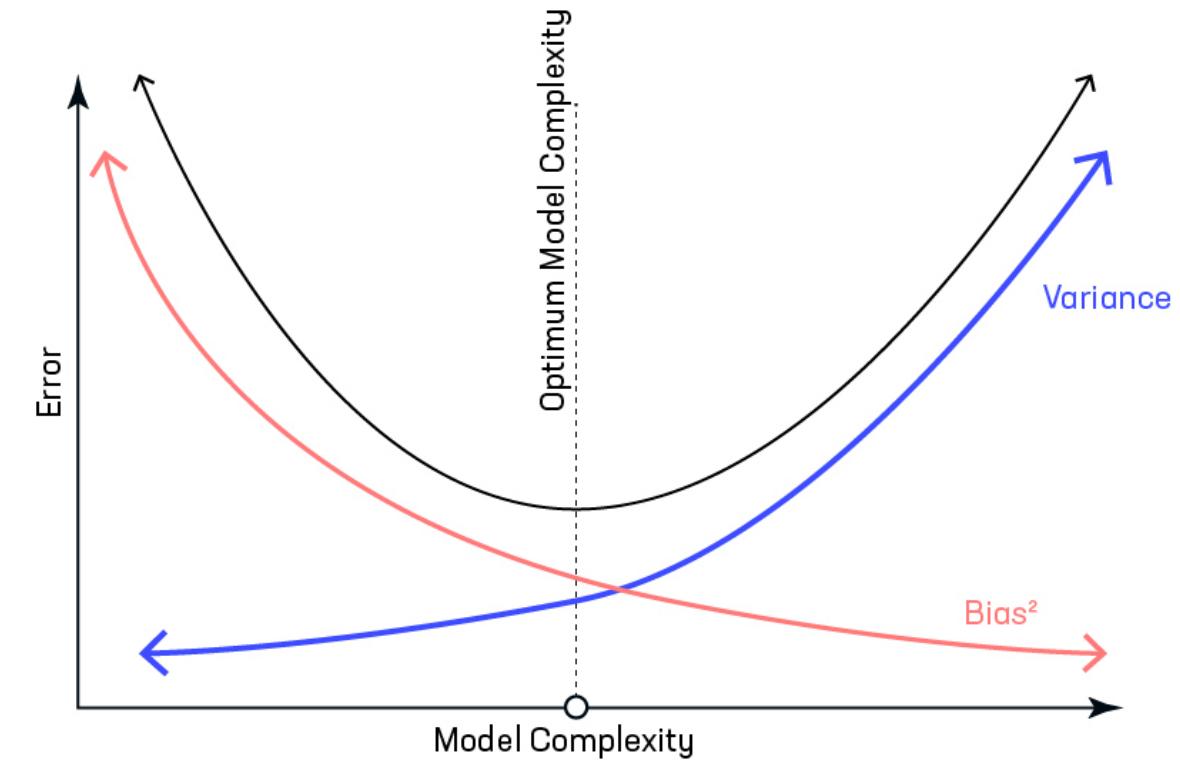
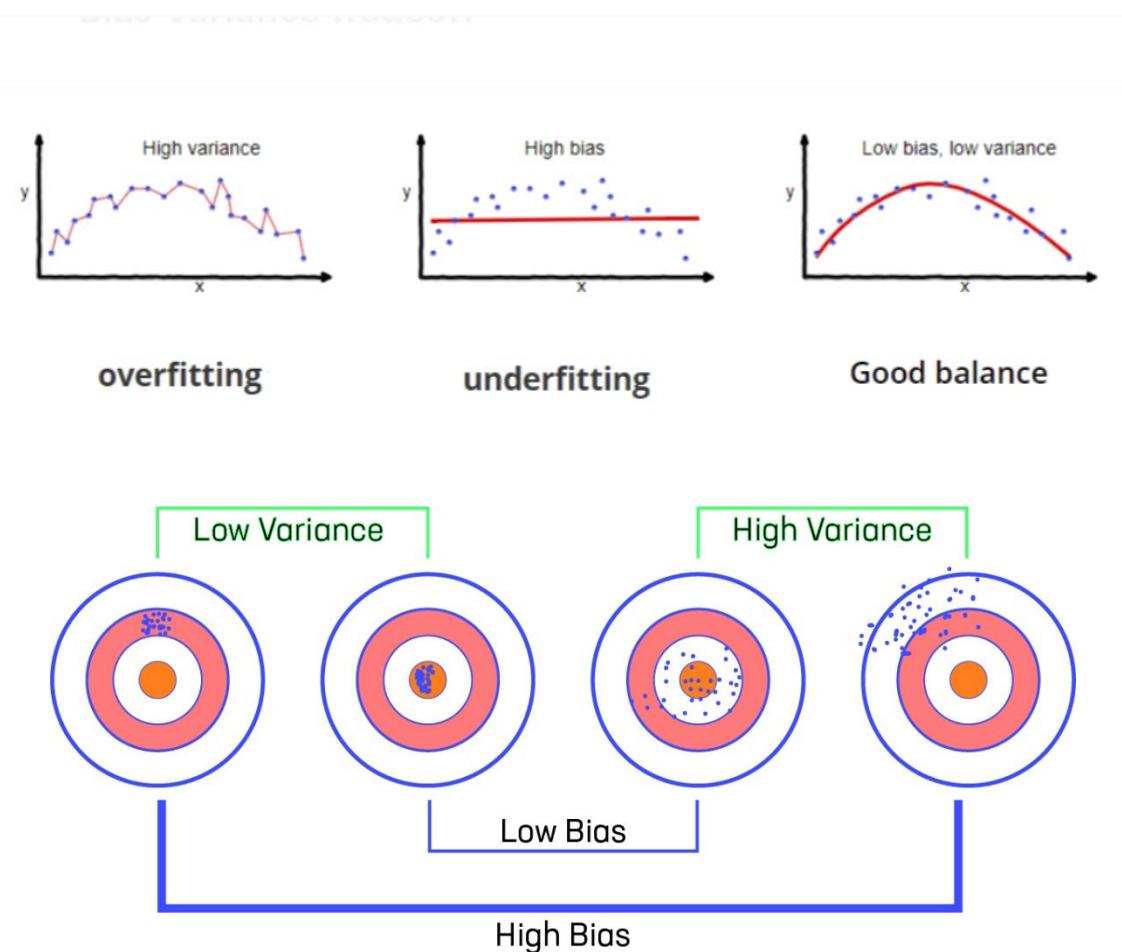


Once again, examine the bias to see if there's a systemic error



The residual from the rational function is improved and less biased
Curve fit parameters determined by Scipy.optimize curve_fit() method

Regularization can and should be used to prevent overfitting





Ridge, Lasso, and Elastic Net are all good regularization options

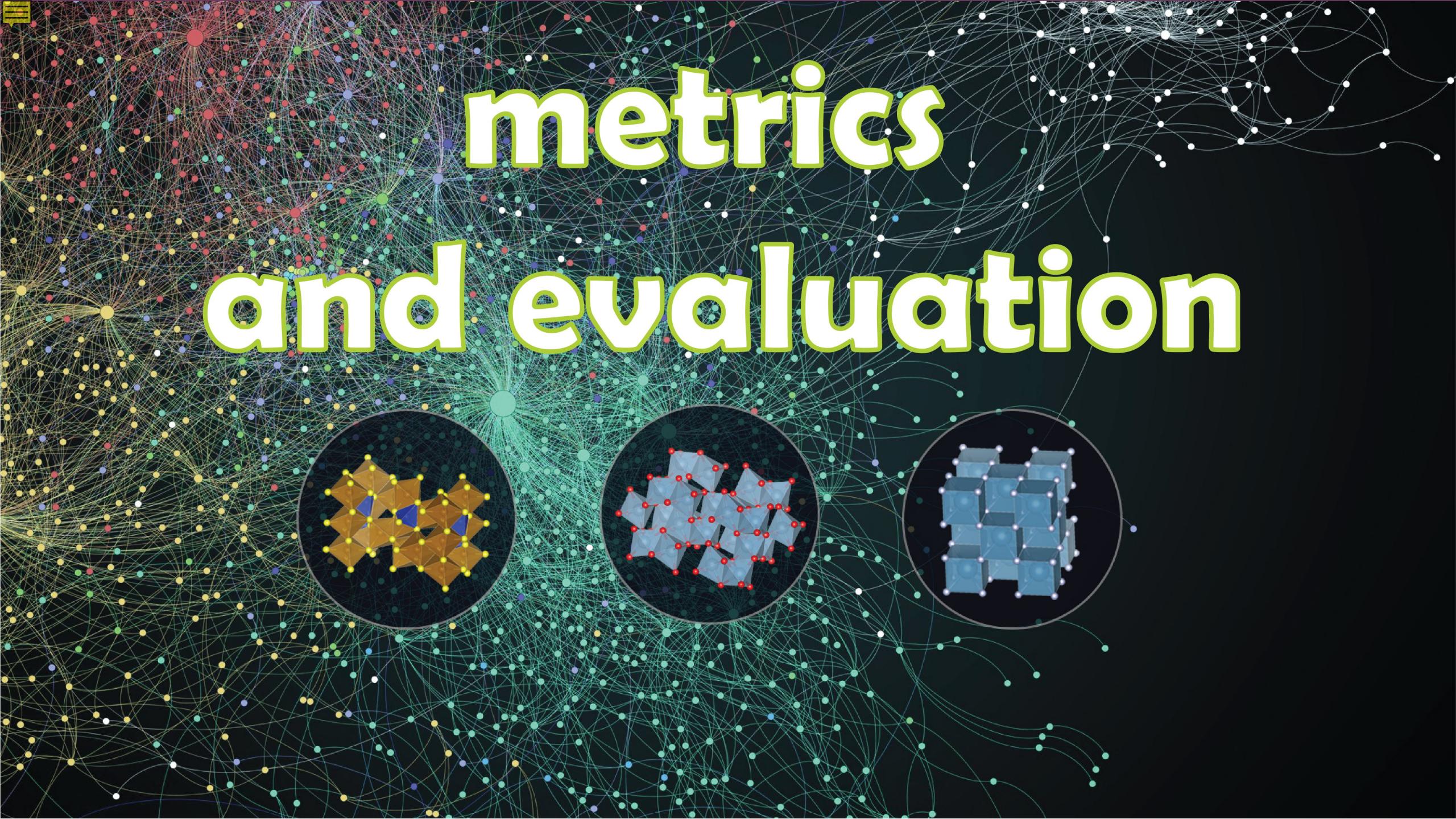
Ridge regression introduces a penalty term based on the square of the coefficients plus a coefficient to control the penalty term. AKA L2 regularization

$$L_{ridge} = \operatorname{argmin}_{\beta} (\|Y - \beta * X\|^2 + \lambda \|\beta\|_2^2)$$

Lasso regression introduces a penalty term based on the sum of the coefficients plus a coefficient to control the penalty term. AKA L1 regularization

$$L_{lasso} = \operatorname{argmin}_{\beta} (\|Y - \beta * X\|^2 + \lambda \|\beta\|_1)$$

Elastic Net combines both.



metrics and evaluation

