

Overfitting

Mohammad Emtiyaz Khan
EPFL

Oct 1, 2015



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

©Mohammad Emtiyaz Khan 2015

Motivation

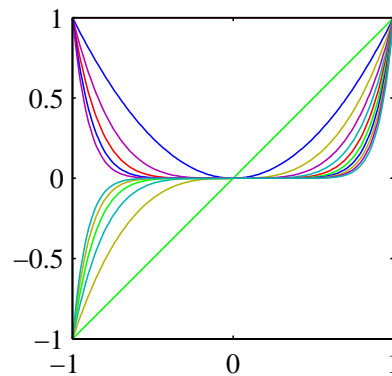
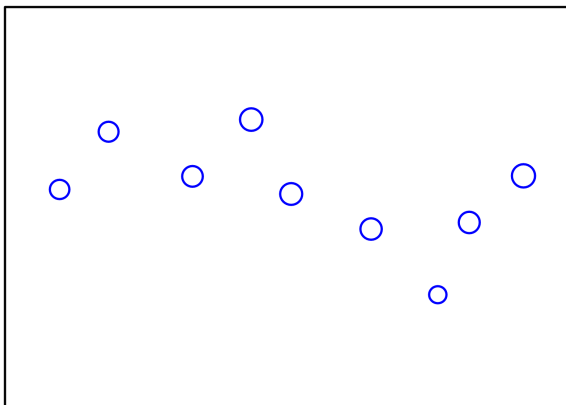
Linear model can be too limited and usually underfit. One way is to use nonlinear basis functions instead.

Consider simple linear regression. Given one-dimensional input x_n , we can generate a polynomial basis.

$$\phi(x_n) = [1, x_n, x_n^2, x_n^3, \dots, x_n^M]$$

Then we fit a linear model using the original *and* the generated features:

$$y_n \approx \beta_0 + \beta_1 x_n + \beta_2 x_n^2 + \dots + \beta_M x_n^M$$

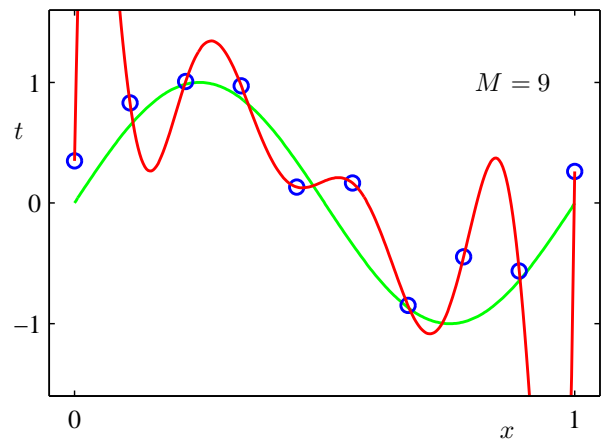
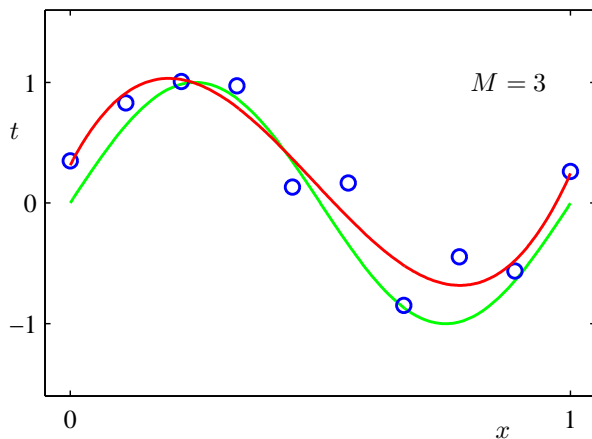
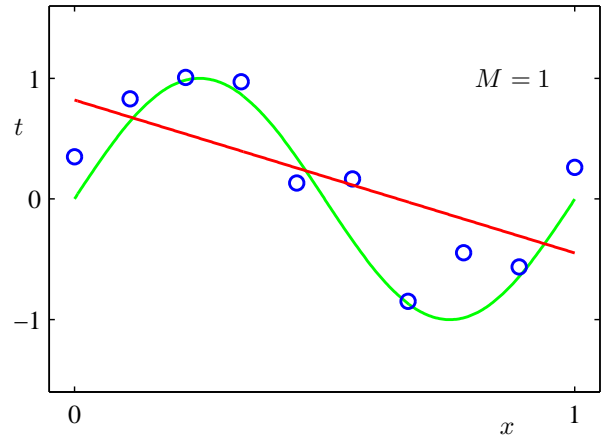
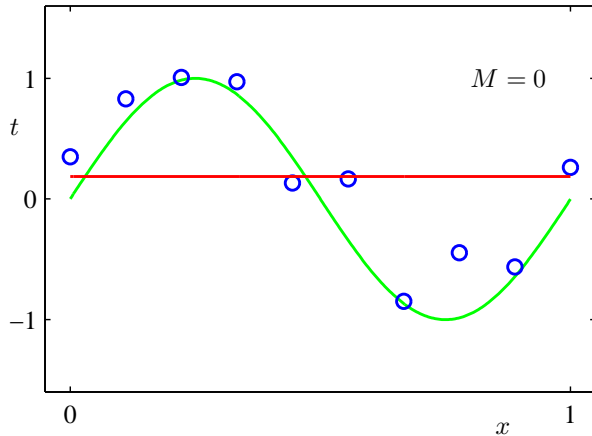


Overfitting and Underfitting

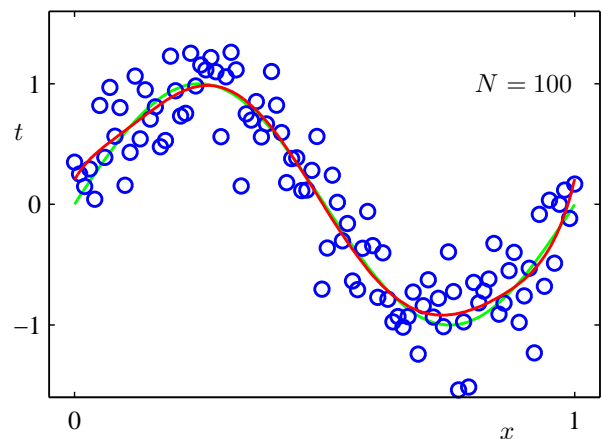
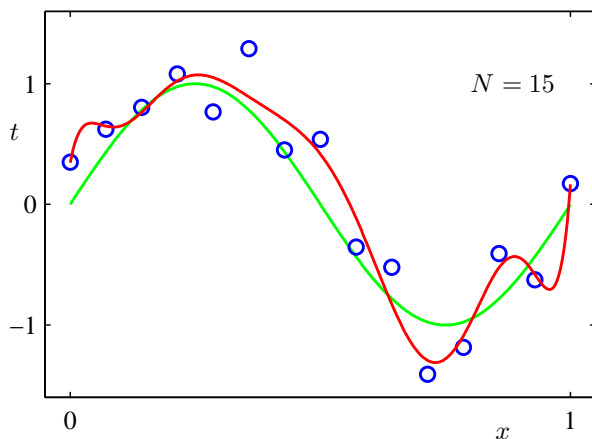
Overfitting is fitting the noise in addition to the signal. **Underfitting** is not fitting the signal well. In reality, it is very difficult to be able to tell the signal from the noise.

Complex models overfit easily

Circles are data points, green line is the truth & red line is the model fit. M is the maximum degree in the generated polynomial basis.



If you increase the amount of data, overfitting *might* reduce.



Occam's razor

One solution is dictated by [Occam's razor](#) which states that “Simpler models are better – in absence of certainty.”

Sometimes, if you increase the amount of data, you might reduce overfitting. But, when unsure, choose a simple model over a complicated one.

We can choose simpler models by adding a [regularization term](#) which ‘penalizes’ complex models.

$$\min_{\beta} \frac{1}{2N} \sum_{n=1}^N [y_n - \tilde{\phi}(\mathbf{x}_n)^T \beta]^2 + \frac{\lambda}{2N} \sum_{j=1}^M \beta_j^2$$

where $\lambda > 0$.

To do

Read about overfitting in the paper by Pedro Domingos (section 3 and 5 of “A few useful things to know about machine learning”).