



Introduction to seminar course

The unreasonable effectiveness of overparameterized machine learning models (3 hp)

Dave Zachariah

Machine Learning Arena IT-CIM & Div. Systems and Control

ML journal clubs



1. Introduction to Machine Learning (2017)
2. Learning Complex Models via Approximate Inference (2018)
3. Adversarial Machine Learning (2019)

ML journal clubs



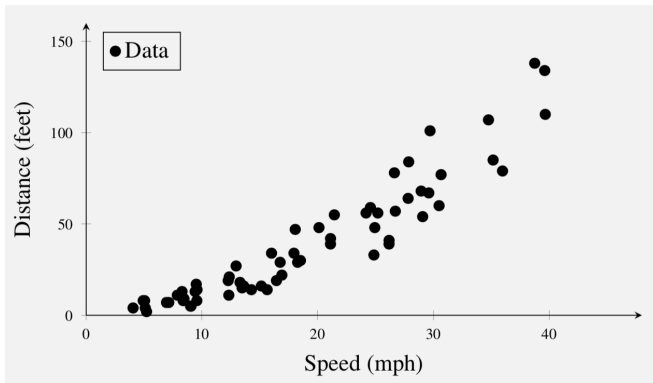
1. Introduction to Machine Learning (2017)
2. Learning Complex Models via Approximate Inference (2018)
3. Adversarial Machine Learning (2019)

New effort: Seminar course!



Background

Textbook example of prediction

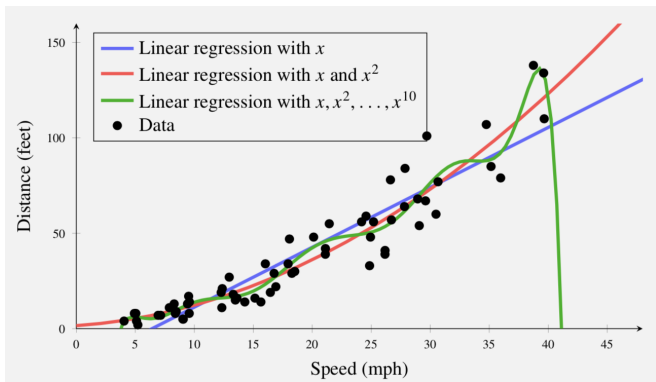


Data on outcome y and covariate x



Lindholm, A. et al.: *Machine Learning - A First Course for Engineers and Scientists*, Cambridge University Press, 2021.

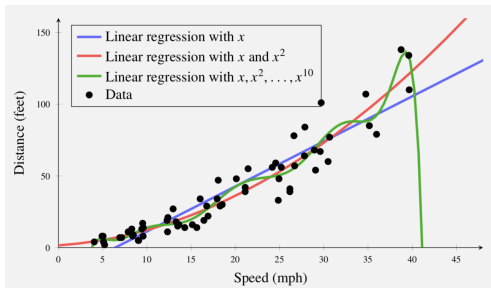
Textbook example of prediction



Example of predictor model parameterized by θ :

$$\hat{y}_{\theta}(x) = \theta^{\top} \begin{bmatrix} 1 \\ x \end{bmatrix} \quad \dim(\theta) = 2$$

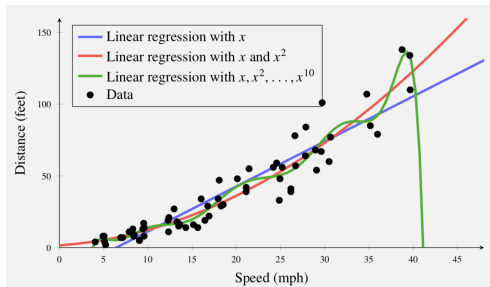
Increasing ‘model complexity’



Example of predictor model parameterized by θ :

$$\hat{y}_{\theta}(x) = \theta^{\top} \begin{bmatrix} 1 \\ x \\ \vdots \\ x^{10} \end{bmatrix} \quad \dim(\theta) = 11$$

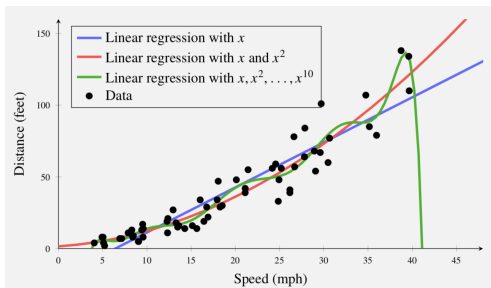
Increasing 'model complexity'



If $\dim(\theta) \leq \# \text{training samples}$, then (typically) there exists **unique** predictor $\hat{y}_{\theta}(x)$ that minimizes training error, e.g.,

$$E_{\text{train}} = \mathbb{E}[(y - \hat{y}_{\theta}(x))^2]$$

Increasing 'model complexity'



If $\dim(\theta) > \# \text{training samples}$, then many predictors $\hat{y}_{\theta}(x)$ can reduce the training error to **zero**, e.g.,

$$E_{\text{train}} = \hat{\mathbb{E}}[(y - \hat{y}_{\theta}(x))^2] = 0$$

Overfit and overparameterization

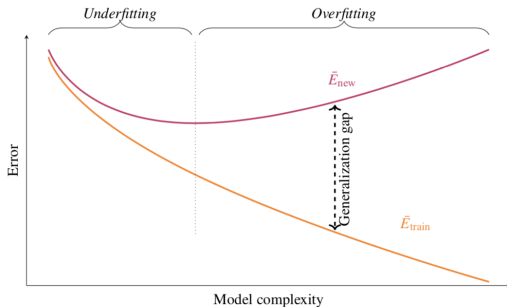


Figure 4.3: Behavior of \bar{E}_{train} and \bar{E}_{new} for many supervised machine learning methods, as a function of model complexity. We have not made a formal definition of complexity, but a rough proxy is the number of parameters that are learned from the data. The difference between the two curves is the generalization gap. The training error \bar{E}_{train}

Training error vs. test error

Overfit and overparameterization

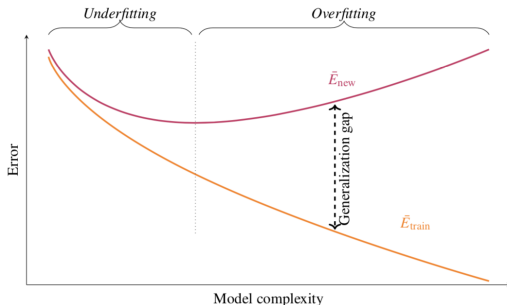


Figure 4.3: Behavior of \bar{E}_{train} and \bar{E}_{new} for many supervised machine learning methods, as a function of model complexity. We have not made a formal definition of complexity, but a rough proxy is the number of parameters that are learned from the data. The difference between the two curves is the generalization gap. The training error \bar{E}_{train}

Puzzle: Modern ML methods using

- ▶ models with $\dim(\theta) \gg \#\text{training samples}$
- ▶ stochastic gradient search methods

achieve **state-of-the-art** test error E_{test} !