

Machine learning algorithms

Overfitting & underfitting

2020-09-29

CSCI 471 / 571, Fall 2020

Kameron Decker Harris

What is ML?

d^p polynomial
Kernel

- Data + Optimization + Statistics → Predictions

\underline{X}, \vec{y}

least squares
normal eqns
- linear system
- solution $\vec{\beta}^*$

fitting

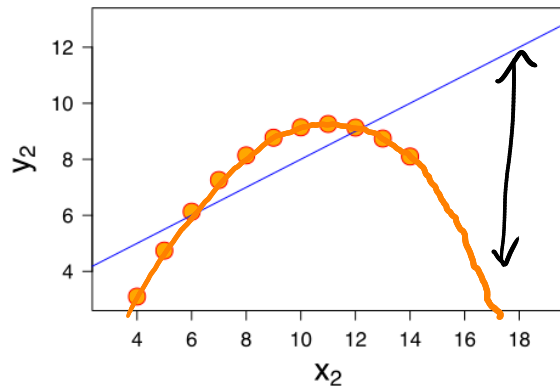
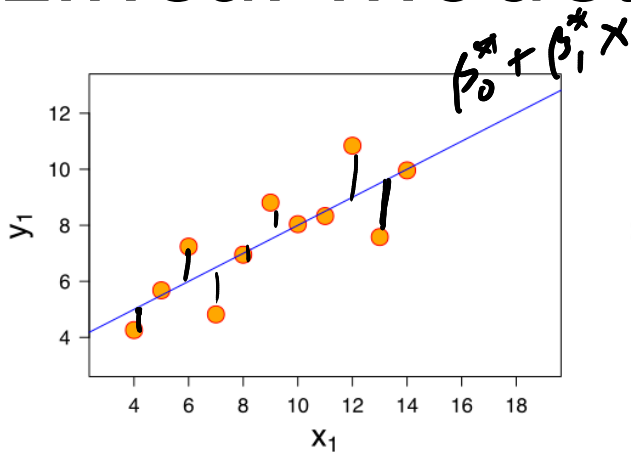
overfitting/under
effects of noise

★ Ref 2. Chapter 3, James, Witten et al.

Linear models are simple

• fits are very close

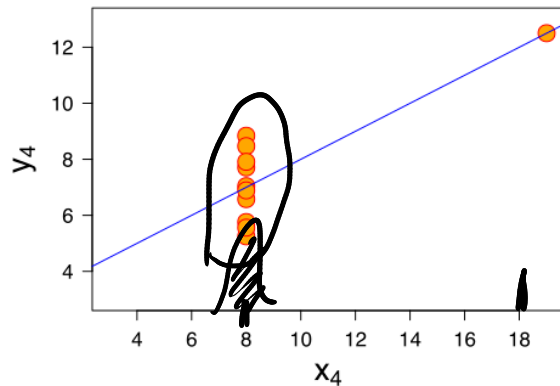
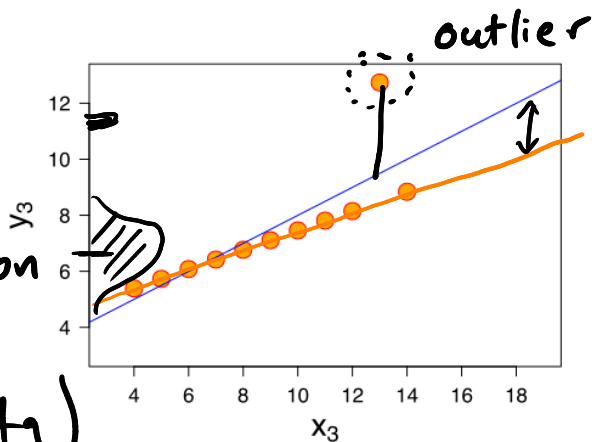
Linear looks good →



★ Some hope for squared err. looks like polynomial

- same sq. err.
- same means \bar{x} \bar{y}
- same std's

poor generalization (predictions outside training data)



Anscombe, 1973

Least-squares in nonlinear basis

$\mathcal{X} \rightarrow$ add polynomial features

$$\min_{\vec{\beta}} \left\| \underline{\Phi}(\mathcal{X}) \vec{\beta} - \vec{y} \right\|^2$$

$n \times p$

e.g.

$$\vec{\Phi}(\vec{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^3 \\ \cos x_2 \end{bmatrix}$$

$d=2$

$$f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^3 + \beta_4 \cos x_2$$

Beware extrapolation

STAT Topics Coronavirus Opinion Podcast Newsletters Reports Events Q

HEALTH

Influential Covid-19 model uses flawed methods and shouldn't guide U.S. policies, critics say

By SHARON BEGLEY @sxbegle / APRIL 17, 2020

[Reprints](#)



