

Last Time:

- Polynomial Regression
- Overfitting vs. underfitting

Ref

ISL: 2.2.1, 2.2.2

ESL: 7.2, 7.3

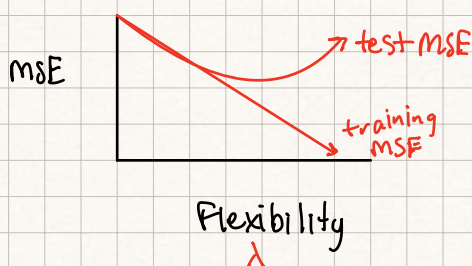
Bishop: 3.2

- Machine Learning finds patterns to data that human intuition cannot find

- Test data is not used in training and is used to measure fit

- Measuring quality of fit

$$MSE = \frac{1}{n} \sum (y_i - f(x_i; w))^2$$



Demo Notes:

overfitting the data results in an unintuitive fit that doesn't accurately model data

Extremely large weights can also indicate overfitting

Regularization: Avoiding Overfitting

Ref: Bishop 3.1.4

LS method: $\min_w \sum (y_i - f(x_i; w))^2$ to find w given $(x_i, y_i), i=1 \dots n$

$$\min_w \left[\sum (y_i - f(x_i; w))^2 + \lambda \sum_{i=1}^n w_i^2 \right] = \min \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

positive regularization constant

regularization term

Also known as hyperparameter

note:

$$\|a\|_2^2 + \|b\|_2^2 = \left\| \begin{bmatrix} a \\ b \end{bmatrix} \right\|_2^2$$

$$= \min_w \left\| \begin{bmatrix} y \\ 0 \end{bmatrix} - \begin{bmatrix} X \\ \sqrt{\lambda} I \end{bmatrix} w \right\|_2^2 \rightarrow \text{standard LS problem}$$

$$w = \left([X^T \sqrt{\lambda} I] \begin{bmatrix} X \\ -\sqrt{\lambda} I \end{bmatrix} \right)^{-1} [X^T \sqrt{\lambda} I] \begin{bmatrix} y \\ 0 \end{bmatrix} = \boxed{(X^T X + \lambda I)^{-1} X^T y}$$

not used numerically

Demo Notes:

weights are reasonable at high order, more intuitive fit

- Ridge Reg.
- Plotting with λ as ind. var.
- As λ inc. MSE inc.
- As λ dec., MSE dec.

$$\text{Lasso: } \min_w \left[\sum (y_i - f(x_i; w))^2 + \lambda \sum |w_i| \right]$$

Finding w is much more complicated

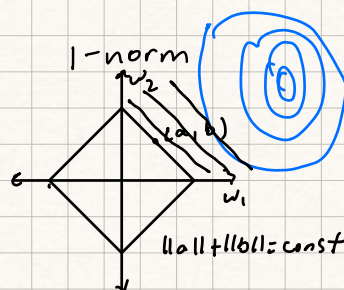
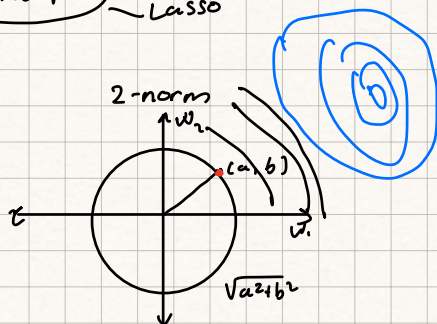
- coordinate descent

the weights are constrained to be small and also sparse

(many zero weights)

Why is solution sparse?

$\lambda \|w\|$ — Lasso



many points
of non-intersection