

## Tutorial Linear Regression

Thursday, 19 November 2020 17:33

**Problem 1:** Assume that we are given a dataset, where each sample  $x_i$  and regression target  $y_i$  is generated according to the following process

$$x_i \sim \text{Uniform}(-10, 10)$$

$$y_i = ax_i^3 + bx_i^2 + cx_i + d + \epsilon_i, \quad \text{where } \epsilon_i \sim \mathcal{N}(0, 1) \text{ and } a, b, c, d \in \mathbb{R}.$$

The 3 regression algorithms below are applied to the given data. Your task is to say what the bias and variance of these models are (low or high). Provide a 1-2 sentence explanation to each of your answers.

a) Linear regression

$d=1$

Bias: high. Variance: low.

A straight line cannot capture a degree 3 polynomial (thus underfitting the data).

b) Polynomial regression with degree 3

$d=3$

Bias: low. Variance: low.

The model is same as the data generating process. We can achieve a good fit.

c) Polynomial regression with degree 10

$d=10$

Bias: low. Variance: high.

Since we are using a polynomial regression with a degree much higher compared to the data generating process, the model will overfit the data.

**Problem 2:** Given is a training set consisting of samples  $X = \{x_1, x_2, \dots, x_N\}$  with respective regression targets  $y = \{y_1, y_2, \dots, y_N\}$  where  $x_i \in \mathbb{R}^D$  and  $y_i \in \mathbb{R}$ .

Alice fits a linear regression model  $f(x_i) = w^T x_i$  to the dataset using the closed form solution for linear regression (normal equations).

Bob has heard that by transforming the inputs  $x_i$  with a vector-valued function  $\Phi$ , he can fit an alternative function,  $g(x_i) = v^T \Phi(x_i)$ , using the same procedure (solving the normal equations). He decides to use a linear transformation  $\Phi(x_i) = A^T x_i$ , where  $A \in \mathbb{R}^{D \times D}$  has full rank.

a) Show that Bob's procedure will fit the same function as Alice's original procedure, that is  $f(x) = g(x)$  for all  $x \in \mathbb{R}^D$  (given that  $w$  and  $v$  minimize the training set error).

Alice uses the normal equation directly and obtains

$$w^* = (X^T X)^{-1} X^T y.$$

Bob fits the model to the transformed data and obtains

$$\begin{aligned} v^* &= ((XA)^T (XA))^{-1} (XA)^T y \\ &= (A^T X^T X A)^{-1} A^T X^T y \quad (\dagger) \\ &= A^{-1} (X^T X)^{-1} (A^T)^{-1} A^T X^T y \quad (\ddagger) \\ &= A^{-1} (X^T X)^{-1} X^T y = A^{-1} w^*. \end{aligned}$$

(Note that  $\Phi$  transforms the column vectors  $x_i$  via  $A^T x_i$  but  $X$  contains the features as rows. Therefore the transformed feature matrix is  $XA$ )

Now it is immediate to see that

$$g(x_i) = v^{*T} \Phi(x_i) = w^{*T} A^{-T} A x_i = w^{*T} x_i = f(x_i).$$

b) Can Bob's procedure lead to a lower training set error than Alice's if the matrix  $A$  is not invertible? Explain your answer.

Any weights  $v^*$  Bob finds are also feasible for Alice by letting  $w = Av^*$ . Therefore Bob can only access a subset of the parameter space and cannot achieve a lower loss value than Alice. (It could still be equal but it cannot be better.)

Note that we are only talking about training error in this example, not test error. Bob might manage to find a model that generalizes better than Alice's, but Alice will always be able to fit the training data at least as well as Bob.

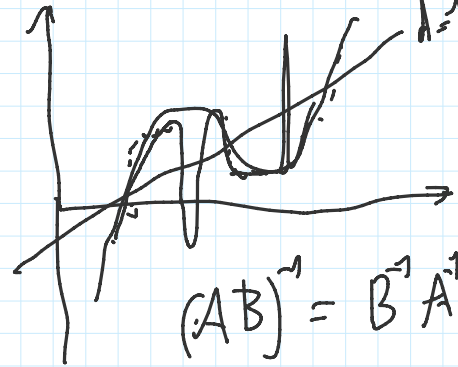
$$x_1, \dots, x_n \quad y_1, \dots, y_n$$

find  $f$  s.t.  $f(x_i) \approx y_i$

and  $f$  generalizes to new  $x$

$$f(x) = w^T \Phi(x)$$

$$\Phi(x) = \begin{pmatrix} x^0 \\ x^1 \\ x^2 \\ \vdots \\ x^d \end{pmatrix}$$



$$(AB)^T = B^T A^T$$

$$E_D = \frac{1}{2} \sum_i (y_i - f(x_i))^2$$

$$= (Xw - y)^T (Xw - y)^{\frac{1}{2}}$$

$$w^* = X^T y = (X^T X)^{-1} X^T y$$

$$\Phi = \begin{pmatrix} x_1^T A \\ \vdots \\ x_n^T A \end{pmatrix}$$

$$v^* = (XA)^T y$$

$$w = A v^*$$

$$g(x) = v^{*T} A^T x$$

$$f(x) = w^T x = v^{*T} A^T x$$