

Lecture 2 Jan 21 th.

What can we do with regression?

- prediction
- fitting

What is the general process of building/training a ml algorithm?

- collect data (care about bias and sample size)

Non-linear Regression //

Basis Function Regression //

Regression Function: $y = f(x) = \sum_k w_k b_k(x)$

where w_k is a weight or the parameters

$b_k(x)$ is a "simple" func. or the basis funcs.

y is the approximation of the relationship

Ex. The regre. function of linear regression is $y = f(x) = \sum_k w_k x^k$

The regre. function of polynomial regression of order 4 is:

$$y = f(x) = w_4 x^4 + w_3 x^3 + w_2 x^2 + w_1 x + w_0$$

Hence the weights are $w_4 \dots w_0$, the basis funcs are $x^4 \dots x^0$

We can rewrite above as

$$\begin{aligned} y = f(x) &= \vec{b}(x)^T \vec{w} \\ &= [b_1(x) \ b_2(x) \ \dots \ b_n(x)] \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \end{aligned}$$

Radial Basis Functions (RBFs) //

A B.F.R. with basis function

$$b_i(x) = \exp\left(-\frac{(x - c_i)^2}{2\sigma^2}\right)$$

$$\begin{aligned} \text{Thus } f(x) &= \sum_k w_k b_k(x) \\ &= \sum_k w_k \exp\left(-\frac{(x - c_k)^2}{2\sigma^2}\right) \end{aligned}$$

Where c_k is the center and σ^2 determines the width of the basis function. (These must be estimated)

To pick the centers, there are 3 methods

- ① Place centers uniformly throughout the data. (can cause impractical number of centers in high dimensions)
- ② Place centers for each data point. (This is used more often since there's a limit for # of centers. (can be expensive for large sample size))
- ③ Cluster the data and place a center for each cluster.

To pick the width,

- ① Trial and error
- ② Use the average square distances to neighbouring centers, scaled by a constant. This allows you to use different widths for each basis function.

Calculation of Error of B.F.R //

We find the error using the least squares method.

$$E(w) = \sum_i (y_i - f(x_i))^2 = \sum_i (y_i - \sum_k w_k b_k(x_i))^2 \\ = \|\vec{y} - B\vec{w}\|^2, \quad \vec{y} \in \mathbb{R}^n, \vec{w} \in \mathbb{R}^m, B \in \mathbb{R}^{n \times m}$$

Which is the same form as $\|\vec{y} - X\vec{w}\|^2$ which is the same form as linear regression. Thus,

$$\vec{w}^* = B^{-1}\vec{y}$$

Issues with Overfitting //

Directly minimizing the square error can lead to overfitting; when the model fits too well with data and can not generalize.

When does this occur:

- ① Problem not sufficiently constrained
- ② More parameters than data
- ③ Fitting noise
- ④ Discarding uncertainty.

There are 2 ways to fix overfitting:

- ① Add prior knowledge
- ② Handling uncertainty

Adding Prior Knowledge : Regularization

Ex: higher degree polynomial always wins out on cross-validation since it can always fit the data the best.

To avoid this, we can regularize the objective function.

$$E(w) = \underbrace{\|\vec{y} - B\vec{w}\|^2}_{\text{data term}} + \lambda \underbrace{\|\vec{w}\|^2}_{\text{smoothness term}}$$

$\|\vec{w}\|$ called weight decay

$$\begin{aligned} &= (\vec{y} - B\vec{w})^T (\vec{y} - B\vec{w}) + \lambda \vec{w}^T \vec{w} \\ &= \vec{w}^T (B^T B + \lambda I) \vec{w} - 2 \vec{w}^T B^T \vec{y} + \vec{y}^T \vec{y} \end{aligned}$$

$$\text{Thus } \vec{w}^* = (B^T B + \lambda I)^{-1} B^T \vec{y}$$

Now the model will become more smooth but will not have the smallest variance.

Artificial Neural Networks.

The basis function chosen is

$$g(a) = \frac{1}{1 + e^{-a}} \quad (\text{Sigmoid Function})$$

The trick here is that

$$\vec{y} = f(\vec{x}) = \sum_j \vec{w}_j g\left(\sum_k w_{kj}^{(1)} x_k + b_j\right) + \vec{b}$$

A linear regression with weights W ($w_{k,j}$) are applied to \vec{x} and is put through a Sigmoid Function and the output is used as features of another linear regression.

$0, 1$ are different values.

Ex. The 1-D Case

$$y = \sum_j w_j g(w_j x + b_j) + b$$

Knn Regression

This method does not need model, does not need fitting
We find the k closests neighbouring points and their (possibly weighted) average.

Ex. Pixel color correction, temperature prediction.

$$y = \frac{1}{k} \sum_{i \in N_k(x)} y_i$$

Weighted Example

$$y = \frac{\sum_{i \in N_k(x)} w(x_i) y_i}{\sum_{i \in N_k(x)} w(x_i)}$$

$$w(x_i) = e^{-\|x_i - x\|^2 / 2\sigma^2}$$

R.B.F.