

How to determine Overfitting and underfitting

So, far we studied about k-fold and Pcr. in order to determine the best K value. Now, the question arises, how can we be 100% sure that the value of K we determined, is neither overfit nor underfit.

This we will understand, with a Graphical technique.

Accuracy = $\frac{\text{# of correctly classify points}}{\text{total # pts}}$

$$\text{Error} = 1 - \text{Accuracy}$$

We always wants to maximize Accuracy and minimize errors.

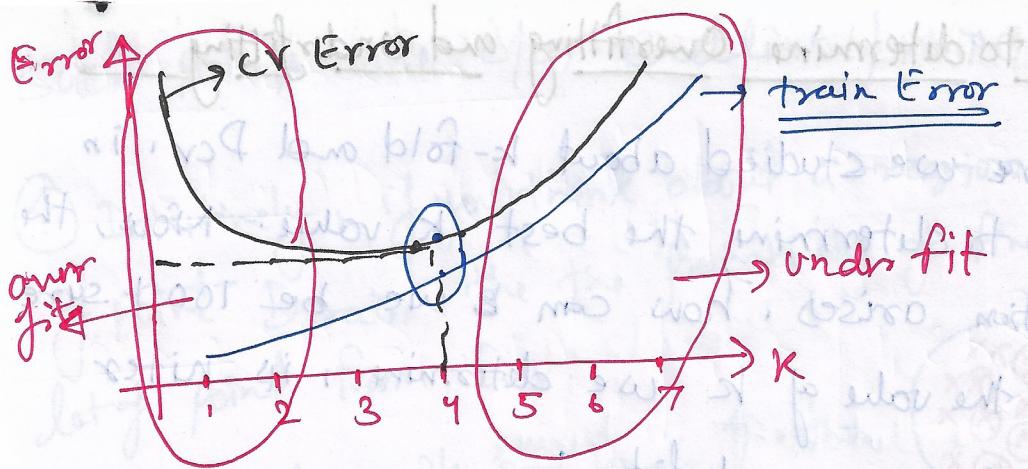
Training Error for each x_i in D_{train} , find

- 2 nearest neighbour to x_i from D_{train} .
- majority vote to get the class label
- if $y_i == \text{classlabel}$ accurate

else error

So, for computing Training Error, we are going to use D_{train} . So, $D_{\text{train}} \rightarrow$

Training Accuracy / Error



As $k \uparrow$ training Error also \uparrow

- If training error is high and CV error is also high then you are under fitting.
- If train error \downarrow and CV error \uparrow then you are overfit.
- For getting the best value of k , or the point where we get the best k , we will have some train error and some cross validation error.

Time based splitting (TBS)

TBS \rightarrow ① Sort your D_n in ascending order of time

② Based on the time, we can split the data oldest 60% $\rightarrow D_{\text{train}}$ \rightarrow first 60%

CV 20% \rightarrow next 20%

test 20% \rightarrow ③ $D_{\text{test}} \rightarrow$ last 20%
 \hookrightarrow latest

One thing that we need to remember is that time is the necessary field for splitting the data in TBS.

We have Amazon fine food review, having time stamp on one column.

Now, let's understand whether TBS or Random splitting is better for this Data set.

Random

$$D_{\text{Train}} \rightarrow \text{NN}$$
$$D_{\text{Cr}} \rightarrow K \text{ in } K-\text{NN}, K=5$$
$$D_{\text{Test}} \rightarrow \text{Accuracy. (93%.)}$$

So, 5-NN accuracy is 93%.

As amazon food review is review of the product. And with time product changes (modified) based on the previous review.

Hence, the new review for a product may also change, so the accuracy ~~of~~ will also change.

So, for Amazon food review Data set, ~~is~~ TBS is better than the Random ~~is~~ select.

Whenever time is available and if things/behaviour/data changes over time, then doing TBS is preferable