

## ★ Cross-validation :

Training set  
↓

This is the portion of dataset used to train the ML model. The model learns patterns and relation from the data in training set.

Validation-set : After training the model on training set, it's important to evaluate its performance on data it hasn't seen before ↓

tuning hyperparameter → i.e.,  $\lambda$  in ridge regression  
learning rate in gradient descent

nos of hidden layer in neural network

Accessing generalization ability - Generalization refer to how well the trained model performed well on unseen data. By evaluating model's performance on validation set we can estimate how well it will generalize to new unseen data

\*\* If model perform well on validation-set it indicate good-generalization ability  
However, if perform poorly it may indicate overfitting (memorizing training data without capturing underlying patterns) or Underfitting (fail to capture underlying pattern in data)-

Test set : reserved for final evaluation of the model's performance

### Types

- \*  $k$ -fold cross validation
- \* Leave-one-out cross validation
- Hold-out validation
- Stratified cross validation

Cross-validation is a technique used in ML to evaluate the performance of a model on unseen data.

It involves dividing the available data into multiple folds or subsets using one of these folds as validation set & training remaining folds.

#### # Leave-one-out cross validation (LOOCV)

Leave- $p$ -out, dataset, Let  $p = 5$

$$n = 100$$

when  $p = 1$

↓  
leave-one-out

iterative check

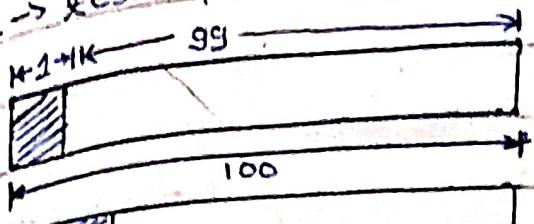
train :  $n-p$

test :  $p$   
(validate-set)

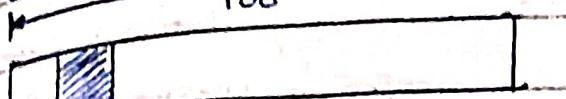
LOOCV is a special case, where nos of folds  $p$  equals to the nos of datapoints in dataset.

The model is trained on  $n-1$  samples, and tested on the one omitted sample, repeating the process for each datapoint in the dataset.

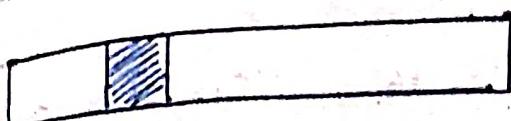
$1 \rightarrow \text{test}, 99 \rightarrow \text{training}$



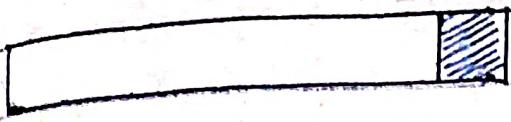
Test 1



Test 2



Test 3



Test 99

- we can make use of all data-points hence it's low bias

- leads to higher variation in the testing model as we are testing against one data-point (outlier - lead to variation)

• lot of execution time

## # K-fold cross validation

We split the dataset into k-numbers of subsets (known as folds) then we perform training on the all the subsets but leave one ( $k-1$ ) subset for the evaluation of the trained model

In this method, we iterate k times with a different subset reserved for testing purpose each time

Total instance : 25

value of k : 5

nos of  
iteration

Training set  
observation

Testing set  
observation

1

[5 6 7 8 9 10 11 12 13 14 15 16]

17 18 19 20 21 22 23 24 25]

[0 1 2 3 4]

2

[0 1 2 3 4 10 11 12 13 14 15 16]

[5 6 7 8 9]

17 18 19 20 21 22 23 24 25]

3

[0 1 2 3 4 5 6 7 8 9 15 16 17]

[10 11 12 13 14]

18 19 20 21 22 23 24 25]

4

[0 1 2 3 4 5 6 7 8 9 15 16 17]

[15 16 17 18 19]

10 18 20 21 22 23 24 25]

5

[0 1 2 3 4 5 6 7 8 9 10 11 12]

[20 21 22 23 24 25]

13 14 15 16 17 18 19]

Note:

It's always suggested that value of k should be 10 as the lower value of k is takes towards validation and higher value of k leads to LOOCV method.

## # Cross-validation

- • overcoming overfitting

• model selection

• hyperparameter tuning

• data-efficient

- • computationally expensive

• time-consuming

• bias-variance tradeoff : The choice of nos of folds in cross-validation can impact b-v tradeoff.

i.e. if too few folds result in high variance

while too many folds lead in high bias

## Cross-Validation Exercise

Q1. Suppose we want to compute 10-fold cross-validation error on 100 training examples. We need to compute error  $N_1$  times, and the cross-validation error is the average of the errors. To compute each error, we need to build a model with data of size  $N_2$ , and test the model on data of size  $N_3$ . What are the appropriate numbers for  $N_1, N_2, N_3$ ?

~~(a)  $N_1 = 10, N_2 = 90, N_3 = 10$~~

$k = 10$

~~(b)  $N_1 = 1, N_2 = 90, N_3 = 10$~~

$N = \text{Samples} = 100$

~~(c)  $N_1 = 10, N_2 = 100, N_3 = 10$~~

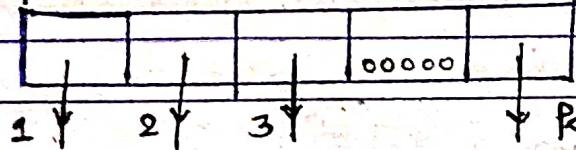
~~(d)  $N_1 = 10, N_2 = 100, N_3 = 100$~~

k-cross validation,  
 $N = \text{samples in training}$

\* generalization error of  
 your model

$N = 100$

\* hyperparameter tuning.



\*  $k$  models are trained  
 and evaluated.

$\frac{N}{k}, \frac{N}{k}, \frac{N}{k}$

model 1 : 1, 2, 3, ...,  $k-1, k$

$k-1$  fold for  
 training

1 fold for  
 testing

model 2 : 1, 2, 3, ...,  $k-1, k$

for testing for training

model  $k$  : 1, 2, 3, ...,  $k-1, k$  → testing

All the folds come in testing at least once

$F_k$  models,  $F_k$  folds

Nos of samples for training

$$\Rightarrow (F_k - 1) \left( \frac{N}{F_k} \right)$$

Nos of samples for testing

$$\Rightarrow (1) \left( \frac{N}{F_k} \right)$$

$$N_1 = ?$$

$N_2$  = training data size

$N_3$  = testing data size

$N_1 \Rightarrow$  Nos of times error is computed = nos of times model trained =  $F_k$

$(F_k - 1)$  folds

total samples in training

$\therefore (F_k - 1) * \text{size of fold}$

$$\Rightarrow (F_k - 1) \times \frac{N}{F_k} \Rightarrow (10 - 1)(100) / 10$$

$$\Rightarrow \frac{9 \times 100}{10} \Rightarrow 90$$

for testing

$$\text{fold } (1) \left( \frac{100}{10} \right) \Rightarrow 10 \text{ samples}$$

$\Rightarrow N_3 -$

(2) Suppose we are performing leave-one-out (LOO) validation and 10-fold cross-validation on a dataset of size 100,000 to pick between 4 different values of a single hyperparameter. How many times greater is the number of models that need to be trained for LOO validation versus 10-fold cross-validation?

Hyperparameter  $\rightarrow$  4 values.

$a_1, a_2, a_3, \dots, a_n$

for every value  $\rightarrow$  train models / CV

compute cross

value with least cross-validation error is chosen

$N = 100,000, \Rightarrow 4 \text{ times}$

How many  $\rightarrow$  models trained  $\rightarrow$  model trained  
times for LOO CV for  $10 = k$  CV

(M1) 10-fold CV model will be trained.  $k = 10 \text{ times}$

$N = 100,000$

10 fold =  $K$

for 4 parameter  $\Rightarrow 4 \times 10$

$\Rightarrow 40 \text{ times}$

(M2) LOO-CV

$N$  samples  $\Rightarrow 1, 2, 3, \dots, N-1, N$   $\overset{\text{test}}{\underset{\text{train}}{\curvearrowleft}}$  model 1

$1, 2, 3, 4, \dots, N-1, N$   $\overset{\text{test}}{\underset{\text{train}}{\curvearrowleft}}$  model 2

$1, 2, 3, \dots, N-1, N$   $\overset{\text{test}}{\underset{\text{train}}{\curvearrowleft}}$  model  $k^{\text{th}}$

testing  $\Rightarrow \frac{N}{K}$  samples here  $\frac{N}{K} = 1$  (100 CV)

$K = N$

No. of models for 100-CV  $\Rightarrow N$

for 4 parameters

$$\rightarrow 4 \times 100,000 \Rightarrow 400,000$$

$$\begin{array}{rcl} \text{100 CV} & \rightarrow & 400,000 \\ \text{k-fold} & \rightarrow & 40 \end{array}$$

$\Rightarrow 10,000$  times greater

③ cross-validation used for



i) evaluate performance of ML model on unseen data



② select hyperparameters



③ to determine the generalization of ML model



on both test / train  
performance

④ To train multiple ML model on diff.  
datasets

C all the model that we train in CV are  
same model