

[Get started](#)[Open in app](#)[Follow](#)

585K Followers



# Cross-Validation in Machine Learning



Prashant Gupta Jun 5, 2017 · 5 min read

There is always a need to validate the stability of your machine learning model. I mean you just can't fit the model to your training data and hope it would accurately work for the real data it has never seen before. *You need some kind of assurance that your model has got most of the patterns from the data correct, and its not picking up too much on the noise, or in other words its low on bias and variance.*

## Validation

[Get started](#)[Open in app](#)

**validation.** Generally, an error estimation for the model is made after training, better known as evaluation of residuals. In this process, a numerical estimate of the difference in predicted and original responses is done, also called the training error. However, *this only gives us an idea about how well our model does on data used to train it.* Now its possible that the model is underfitting or overfitting the data. So, the **problem with this evaluation technique is that it does not give an indication of how well the learner will generalize to an independent/ unseen data set.** Getting this idea about our model is known as Cross Validation.

## Holdout Method

Now a *basic remedy for this involves removing a part of the training data and using it to get predictions from the model trained on rest of the data.* The error estimation then tells how our model is doing on unseen data or the validation set. **This is a simple kind of cross validation technique, also known as the holdout method.** Although this method doesn't take any overhead to compute and is better than traditional validation, *it still suffers from issues of high variance. This is because it is not certain which data points will end up in the validation set and the result might be entirely different for different sets.*

## K-Fold Cross Validation

As there is never enough data to train your model, *removing a part of it for validation poses a problem of underfitting. By reducing the training data, we risk losing important patterns/ trends in data set, which in turn increases error induced by bias.* So, what we require is a method that provides ample data for training the model and also leaves ample data for validation. K Fold cross validation does exactly that.

In **K Fold cross validation**, the data is divided into  $k$  subsets. Now the holdout method is repeated  $k$  times, such that *each time, one of the  $k$  subsets is used as the test set/ validation set and the other  $k-1$  subsets are put together to form a training set.* The error estimation is averaged over all  $k$  trials to get total effectiveness of our model. As can be seen, every data point gets to be in a validation set exactly once, and gets to be in a training set  $k-1$  times. *This significantly reduces bias as we are using most of the data for fitting, and also significantly reduces variance as most of the data is also being used in validation set.* Interchanging the training and test sets also adds to the effectiveness of this method. **As a general rule and empirical evidence,  $K = 5$  or  $10$  is generally preferred,** but nothing's fixed and it can take any value.

[Get started](#)[Open in app](#)

dataset concerning price of houses, there might be large number of houses having high price. Or in case of classification, there might be several times more negative samples than positive samples. For such problems, *a slight variation in the K Fold cross validation technique is made, such that each fold contains approximately the same percentage of samples of each target class as the complete set, or in case of prediction problems, the mean response value is approximately equal in all the folds.* This variation is also known as **Stratified K Fold**.

*Above explained validation techniques are also referred to as Non-exhaustive cross validation methods.* These do not compute all ways of splitting the original sample, i.e. you just have to decide how many subsets need to be made. Also, these are approximations of *method explained below, also called Exhaustive Methods, that computes all possible ways the data can be split into training and test sets.*

## Leave-P-Out Cross Validation

This approach leaves  $p$  data points out of training data, i.e. if there are  $n$  data points in the original sample then,  $n-p$  samples are used to train the model and  $p$  points are used as the validation set. This is repeated for all combinations in which original sample can be separated this way, and then the error is averaged for all trials, to give overall effectiveness.

*This method is exhaustive in the sense that it needs to train and validate the model for all possible combinations, and for moderately large  $p$ , it can become computationally infeasible.*

**A particular case of this method is when  $p = 1$ . This is known as Leave one out cross validation.** This method is generally preferred over the previous one because *it does not suffer from the intensive computation, as number of possible combinations is equal to number of data points in original sample or  $n$ .*

Cross Validation is *a very useful technique for assessing the effectiveness of your model, particularly in cases where you need to mitigate overfitting.* It is *also of use in determining the hyper parameters of your model*, in the sense that which parameters will result in lowest test error. This is all the basic you need to get started with cross validation. You can get started with all kinds of validation techniques using **Scikit-Learn**, that gets you up and running with just a few lines of code in python.

[Get started](#)[Open in app](#)

If you liked this article, be sure to click ♥ below to recommend it and if you have any questions, **leave a comment** and I will do my best to answer.

For being more aware of the world of machine learning, **follow me**. It's the best way to find out when I write more articles like this.

You can also follow me on [Twitter](#), [email me directly](#) or [find me on linkedin](#). I'd love to hear from you.

That's all folks, Have a nice day :)

---

## Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look](#).

[Get this newsletter](#)

You'll need to sign in or create an account to receive this newsletter.

[Machine Learning](#)[Deep Learning](#)[Artificial Intelligence](#)[Data Science](#)[Data](#)[About](#) [Help](#) [Legal](#)

Get the Medium app

