



# SNS COLLEGE OF TECHNOLOGY

(An Autonomous Institution)

Coimbatore – 35.



## DEPARTMENT OF BIOMEDICAL ENGINEERING

### UNIT 2

#### VC DIMENSION

VC dimension, short for Vapnik-Chervonenkis dimension, is a measure of the complexity of a machine learning model. It is named after the mathematicians Vladimir Vapnik and Alexey Chervonenkis, who developed the concept in the 1970s as part of their work on statistical learning theory.

VC dimension is defined as the largest number of points that can be shattered by a binary classifier without misclassification. In other words, it is a measure of the model's capacity to fit arbitrary labeled datasets. The more complex the model, the higher its VC dimension.

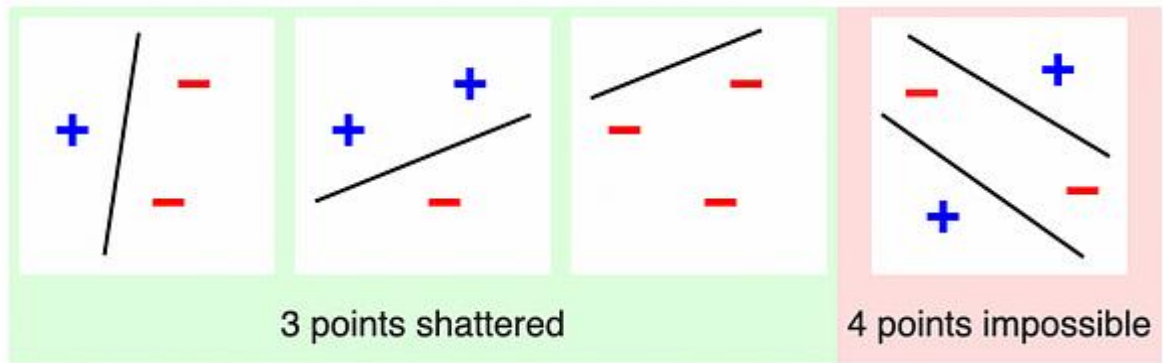
**Mathematically, the VC dimension of a binary classifier is defined as follows:**

Given a set of  $n$  points  $S = \{x_1, x_2, \dots, x_n\}$  in a  $d$ -dimensional space and a binary classifier  $h$ , the VC dimension of  $h$  is the largest integer  $d$  such that there exists a set of  $d$  points that can be shattered by  $h$ , i.e., for any labeling of the  $d$  points, there exists a hypothesis  $h$  in  $H$  that correctly classifies them.

Formally, the VC dimension of  $h$  is:

$$VC(h) = \max\{d \mid \text{there exists a set of } d \text{ points that can be shattered by } h\}$$

*VC dimension has important implications for machine learning models. It is related to the model's generalization ability, i.e., its ability to perform well on unseen data. A model with a low VC dimension is less complex and is more likely to generalize well, while a model with a high VC dimension is more complex and is more likely to overfit the training data.*



Img Src: [https://en.wikipedia.org/wiki/Vapnik%E2%80%93Chervonenkis\\_dimension](https://en.wikipedia.org/wiki/Vapnik%E2%80%93Chervonenkis_dimension)

VC dimension is used in various areas of machine learning, such as support vector machines (SVMs), neural networks, decision trees, and boosting algorithms. In SVMs, the VC dimension is used to bound the generalization error of the model. In neural networks, the VC dimension is related to the number of parameters in the model and is used to determine the optimal number of hidden layers and neurons. In decision trees, the VC dimension is used to measure the complexity of the tree and to prevent overfitting.

#### **Limitations to VC Dimension:**

However, there are some limitations to VC dimension. First, it only applies to binary classifiers and cannot be used for multi-class classification or regression problems. Second, it assumes that the data is linearly separable, which is not always the case in real-world datasets. Third, it does not take into account the distribution of the data and the noise level in the dataset.

The most powerful model according to VC dimension is the SVM with a Gaussian kernel, which has an infinite VC dimension. This means that the model has the capacity to fit any labeled dataset perfectly. However, this does not necessarily mean that it is the best model for all problems, as it can be computationally expensive and may not generalize well to unseen data.

**Finite VC Dimension:**

On the other hand, decision trees have a finite VC dimension, which makes them less complex and more likely to generalize well. However, they may not be suitable for datasets with high dimensionality or complex decision boundaries.

In summary, VC dimension is a powerful tool for measuring the complexity of machine learning models and for understanding their generalization ability. It has important implications for model selection, regularization, and optimization. However, it has some limitations and should be used in conjunction with other evaluation metrics and techniques.

**Reference:**

<https://medium.com/@chandu.bathula16/machine-learning-concept-71-vapnik-chervonenkis-dimension-vc-dimension-211c2c831518>