# Module 3
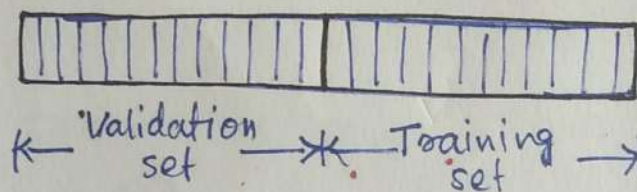
## Cross validation and re-sampling methods

* To test the performance of a classifier, we need to have a number of training/validation set pairs.



← Validation set → ← Training set →

* Cross validation methods are used for generating multiple training - validation sets from a given dataset.

* Cross validation is a technique to evaluate predictive models by partitioning the original sample into a training set to train the model, and a test set to evaluate it.

Different methods are —

→ 1) Hold out method.

→ 2) K-fold cross validation.

→ 3) Leave-one-out cross validation (LOOCV).

→ 4) Bootstrapping.

1) Holdout Method

* Simplest kind of cross validation.

* The dataset is seperated into two sets, called the training set and the testing set.

* The algorithm fits a function using the training set only. Then the function is used to predict the output values for the data in the testing set.

Advantages

→ Simple and easy to run

→ Lower computational cost as it only needs to be run once.

Disadvantages

→ Only work on large dataset.

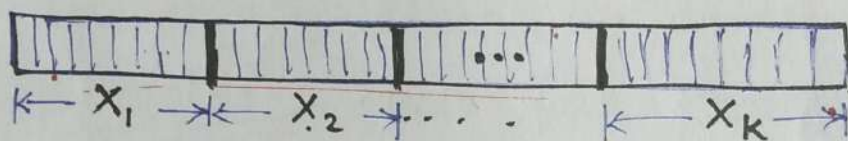→ Higher variance given the smaller size of the data.

2) K-fold cross-validation

* The dataset X is divided randomly into K equal-sized parts, $X_i$ & $i = 1, ---, K$.

* To generate each pair, we keep one of the k parts out as the validation set $V_i$, and combine ~~the remaining~~ the remaining $K-1$ parts to form the training set, $T_i$.

* Doing this K times, we get K pairs $(V_i, T_i)$.



$$\leftarrow X_1 \rightarrow \leftarrow X_2 \rightarrow \cdots \quad \leftarrow X_K \rightarrow$$

1st  $\quad V_1 = X_1 \qquad T_1 = X_2 \cup X_3 \cup \ldots \cup X_K. \qquad P_1 = ?$

2nd  $\quad V_2 = X_2 \qquad T_2 = X_1 \cup X_3 \cup \ldots \cup X_K \qquad P_2 = ?$

$K^{th} \quad V_K = X_K \qquad T_k = X_1 \cup X_2 \ldots \cup X_{K-1} \qquad P_K = ?$

$$\boxed{P = \sum_{i=1}^{k} P_i}$$

Problems with this approach:

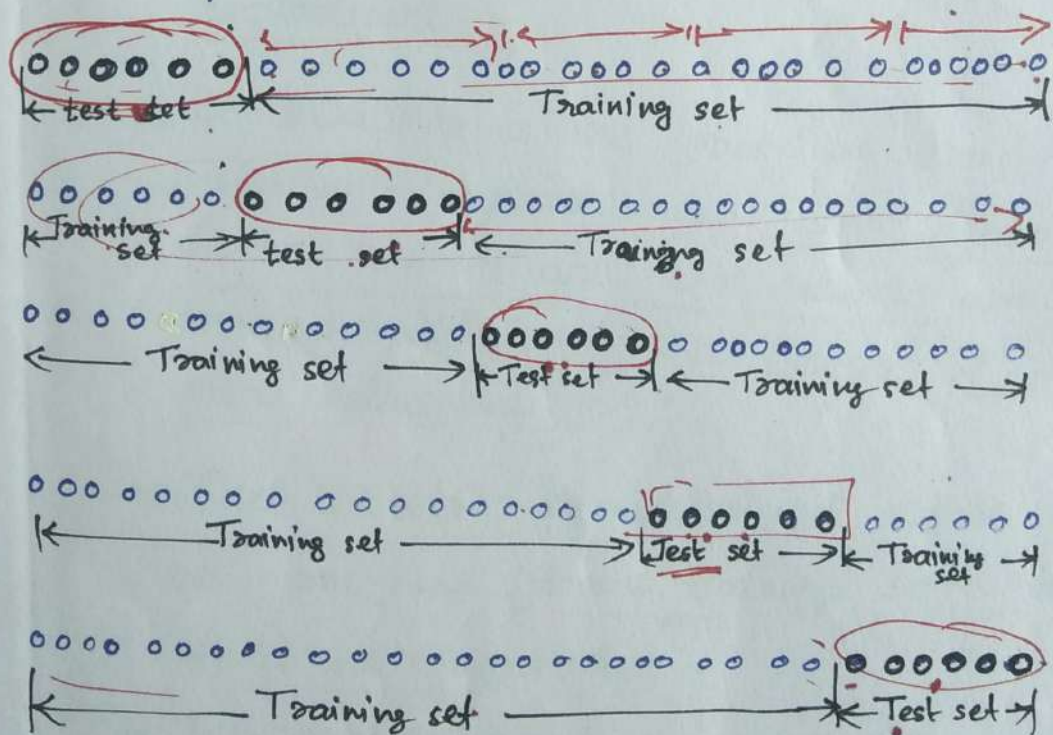→ To keep the training set large, we allow validation sets to be small:

→ Every two training sets share $K-2$ parts.

\* K is typically 10 or 30. As k increases, the percentage of training instances increases. and we get more robust estimators; But the validation set becomes smaller. Also the cost of training the classifier increases as k increases.
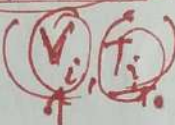
## Example:

Consider a dataset containing 30 samples. And let K = 5. Then we divide dataset into 5 folds, each fold containing 6 samples.

## 3) Leave - one-out Cross validation (LOOCV)

* Given a dataset of N instances, only one instance is left out as the validation set and remaining N-1 instances are used for training.

* We get (N pairs) and hence N iterations are performed. ($V_i, f_i$)

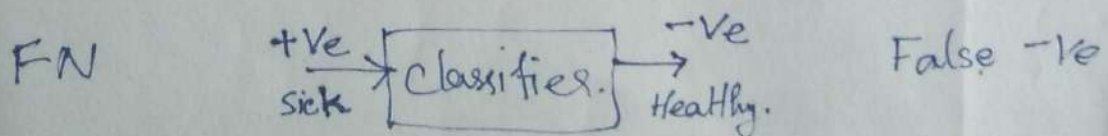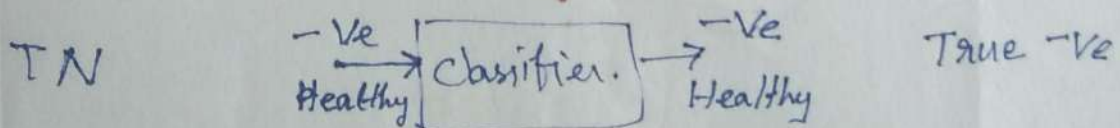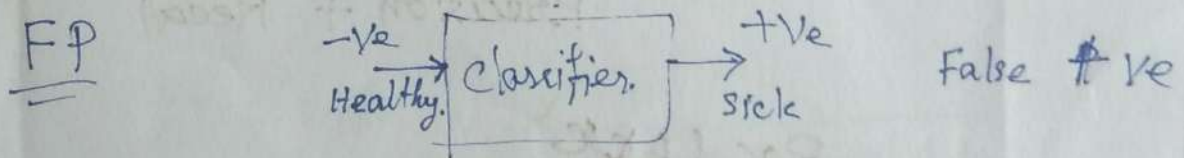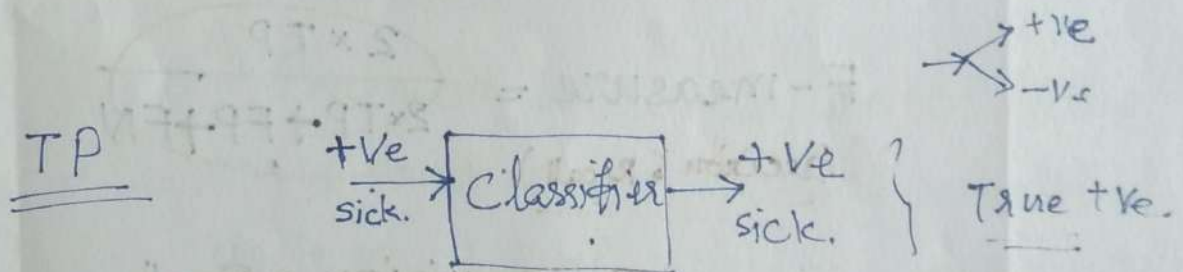## 4) Bootstrapping

* Also known as bootstrap sampling, bootstrap, ~~bo~~ or random sampling with replacement.

* Bootstrapping is the process of computing performance measures using several randomly selected training and test datasets which are selected through a process of sampling with replacement.

* The bootstrap procedure will create one or more new training datasets some of which are repeated.

* The corresponding test datasets are then constructed from the set of examples that were not selected for the respective training datasets.

# Measuring Classifier Performance

Test set → Model → Performance ?

{ True positive (TP)
False positive (FP)
True negative (TN)
False negative (FN) }

→ Precision
→ Recall
→ Accuracy
→ Error rate.
→ Sensitivity
→ Specificity
→ F-measure.
→ ROC & AUC

→ Confusion matrix ?

→ +ve
→ -ve

TP | +Ve sick. → Classifier → +Ve sick. ? | True +Ve.

FP | -Ve Healthy. → Classifier. → +Ve Sick | False +Ve

TN | -Ve Healthy → Classifier. → -Ve Healthy | True -Ve

FN | +Ve sick → Classifier. → -Ve Healthy. | False -Ve

# Confusion Matrix

$(2 \times 2)$ → TP  
FP  
FN  
TN

|  | Actual condition is +ve | Actual condition is −ve |
|---|---|---|
| Predicted condition is +Ve | TP | FP |
| Predicted condition is −Ve | FN | TN |

## Precision & Recall

**Precision** — It is the ratio between true positives and all the predicted positives.

* Precision is a measure of how many samples are correctly identified as +ve out of all the samples which are predicted as +ve.

$$\text{Precision} = \frac{\text{No. of TP}}{\text{Total No. of Predicted +ve}}$$

$$= \frac{TP}{TP + FP}$$

Recall — <u>It</u> is the ratio between true positives and all the actual positives.

* Recall is a measure of how many samples are correctly identified as +ve out of all the samples which are actually +ve.

$$\text{Recall} = \frac{\text{No. of TP}}{\text{Total no. of actual positives}}$$

$$= \frac{TP}{TP + FN}$$

## Other measures :

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Error\ rate = 1 - Accuracy.$$

$$Sensitivity = \frac{TP}{TP + FN} \quad \left\{ \begin{array}{l} Recall \\ TPR \end{array} \right.$$

$$Specificity = \frac{TN}{TN + FP}$$

$$F-measure = \boxed{\frac{2 \times TP}{2 \times TP + FP + FN}}$$
(Precision & Recall)

$$= \frac{2 \times precision \times Recall}{Precision + Recall}$$

$$\boxed{Roc\ \&\ AUC}$$

# — Precision & Recall —
## Problems

## Problem 1

Suppose a computer program for recognizing dogs in photographs identifies eight dogs in a picture containing 12 dogs and some cats. Of the eight dogs identified, five actually are dogs while the rest are cats. Compute the precision and recall of the computer program.

**Ans)**

|  | Actual dogs | Actual cats |  |
|---|---|---|---|
| Predicted dogs | T.P   5 | FP   3 | ⑧ |
| Predicted cats | FN   7 | TN |  |
|  | 12 |  |  |

dog → +ve
cat → −ve

$$\text{precision} = \frac{TP}{\text{No. predicted positive}}$$

$$= \frac{TP}{TP + FP} = \frac{5}{5+3} = \frac{5}{8}$$

$$\text{Recall} = \frac{TP}{\text{Actual Positive}} = \frac{TP}{TP + FN}$$

$$= \frac{5}{5+7} = \frac{5}{12}$$