

11. **F-measure (F_1 or F-score)** : Harmonic mean of precision and recall,

$$F = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

12. F_β : Weighted measure of precision and recall and assigns β times as much weight to recall as to precision.

$$F_\beta = \frac{(1 + \beta^2) \times \text{Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}}$$

where β is a non-negative real number.

13. Classifiers can also be compared with respect to :

- Speed
- Robustness
- Scalability
- Interpretability

14. **Re-substitution error rate**

- Re-substitution error rate is a performance measure and is equivalent to training data error rate.
- It is difficult to get 0% error rate but it can be minimized, so low error rate is always preferable.

Syllabus Topic : Holdout

4.4.2 Holdout

- In holdout method, data is divided into training data set and testing data set (usually 1/3 for testing, 2/3 for training).

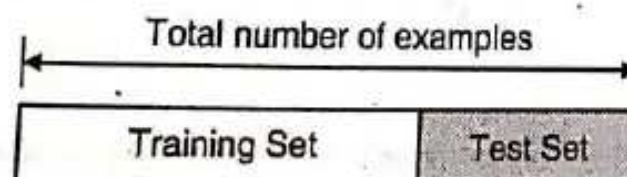


Fig. 4.4.1

to train the classifier, training data set is used and once the classifier is constructed then use test data set to estimate the error rate of the classifier.

Scanned by CamScanner

- If the training is more than better model is constructed and if the test data is more than more accurate the error estimates.
- **Problem :** The samples might not be representative. For example, some classes might be represented with very few instances or even with no instances at all.
- **Solution :** stratification is the method which ensures that both training and testing data have equal number of samples of same class.

Syllabus Topic : Random Sampling

4.4.3 Random Subsampling

- It is a variation of the holdout method.
- The holdout method is repeated k times.
- Each split randomly selects a fixed number example without replacement.

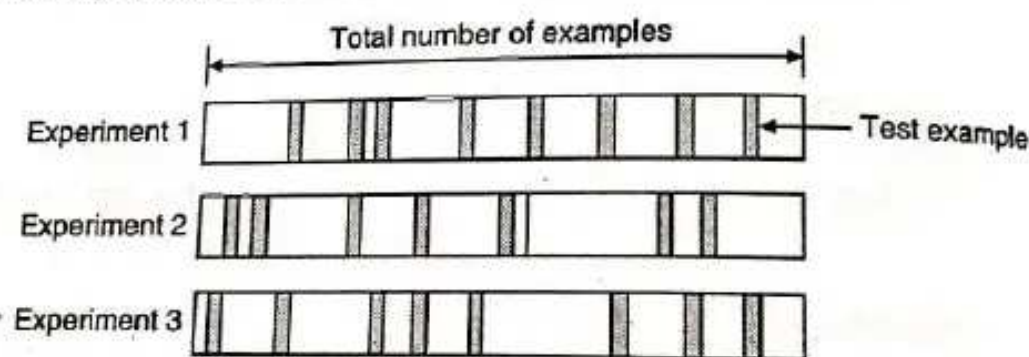


Fig. 4.4.2

- For each data split we retrain the classifier from scratch with the training examples and estimate E_i with the test examples.
- The overall accuracy is calculated by taking the average of the accuracies obtained from each iteration.

$$E = \frac{1}{K} \sum_{i=1}^K E_i$$

Syllabus Topic : Cross-Validation

4.4.4 Cross-Validation (CV)

- Avoids overlapping test sets.

Scanned by CamScanner

k-fold cross-validation

- o First step : Data is split into k subsets of equal size (usually by random sampling).

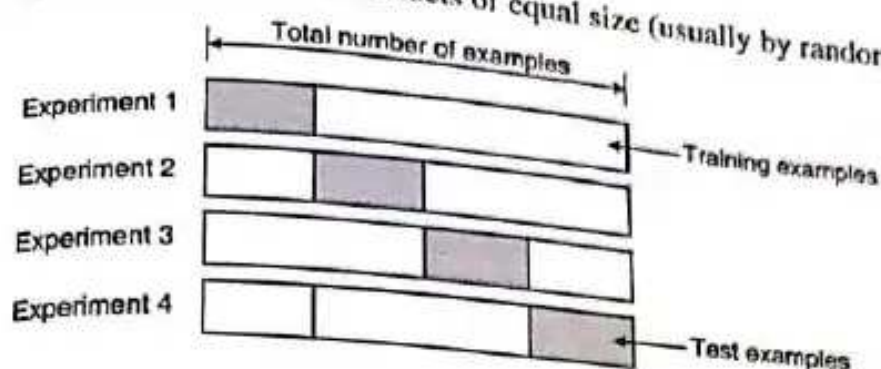


Fig. 4.4.3

- o Second step : Each subset in turn is used for testing and the remainder for training.

The advantage is that all the examples are used for both training and testing.

The error estimates are averaged to yield an overall error estimate.

$$E = \frac{1}{K} \sum_{i=1}^K E_i$$

Leave-one-out cross validation

- o If dataset has N examples, then N experiments to be performed for Leave-one-out cross validation.

- o For every experiment, training uses N-1 examples and remaining example for testing.

- The average error rate on test examples gives the true error.

$$E = \frac{1}{N} \sum_{i=1}^N E_i$$

- **Stratified cross-validation:** Subsets are stratified before the cross-validation is performed.

Stratified ten-fold cross-validation

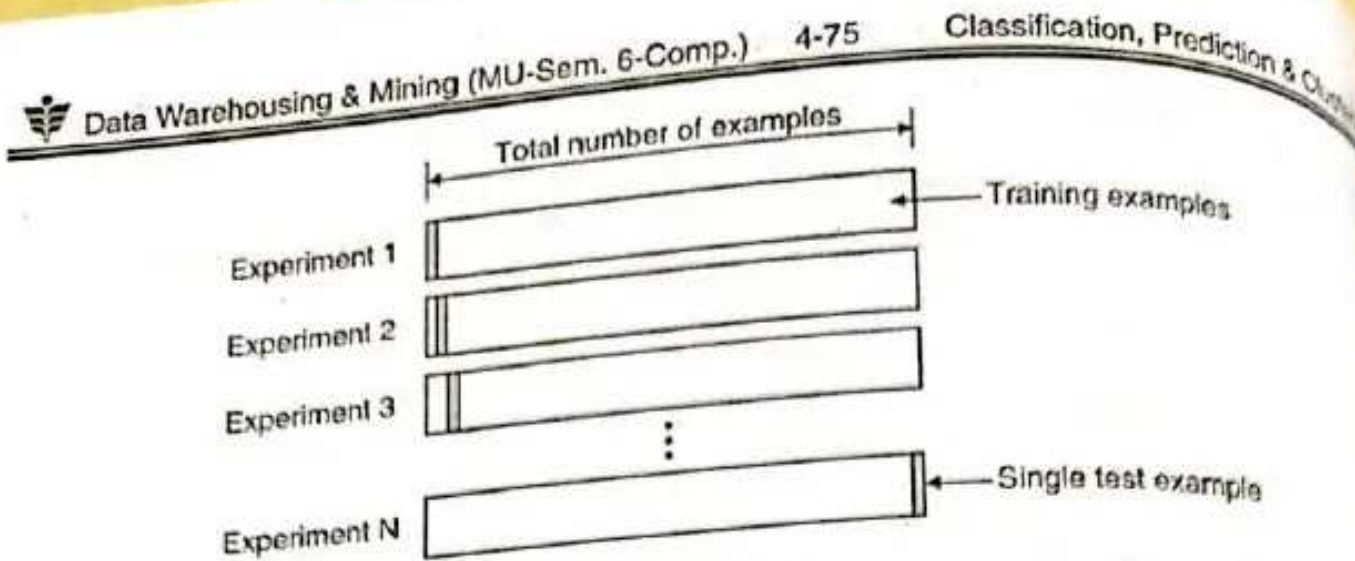
- o This gives accurate estimate of evaluation.

- o The estimate's variance get reduced due to stratification.

- o Ten-fold cross-validation is repeated ten times and finally the results are averaged

based on the previous 10 results.

Scanned by CamScanner



Syllabus Topic : Bootstrap

4.4.5 Bootstrapping

- CV uses sampling of data set without replacement. Once the tuple or instance is selected, it cannot be selected again for training or test data.
- The bootstrap uses sampling with replacement to get the training set.
- **Training set** : A dataset of k instances is sampled with replacement k times to form the training set of k instances.
- **Test set** : This is separate dataset from the original dataset which is not the part of training dataset.
- Bootstrapping is the best error estimator for small datasets.

Syllabus Topic : Clustering

4.5 What is Clustering ?

4.5.1 What is Clustering ?

→ (MU - Dec. 2010, May 2012, Dec. 2012, Dec. 2013)

- Clustering is an unsupervised learning problem.

Data Warehousing & Mining

- A cluster contains data points that are similar to each other.
- We can show this with a diagram.



- From Fig. 4.5 Geometrical data points belong to which cluster.
- The other kind of cluster are a descriptive cluster.

Applications

Clustering is used in

- Marketing : A database, a customer can be identified.
- Biology : A set of classes can be identified.
- Libraries : A set of books can be identified.
- Insurance : A set of policies can be identified.
- City : A set of cities can be identified.