

Machine Learning

Unit 1

* Well Posed problems

→ machine learning is a multidisciplinary field that draws on results from artificial intelligence, probability and statistics, computational complexity theory, control theory, information theory, philosophy, psychology and other fields.

→ Definition of learning - A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .

→ In general, to have a well-defined learning problem, we must identify these three features: the class of tasks, the measure of performance to be improved, and source of experience.

→ Eg checker's problem, robot driving and handwriting recognition

* The machine learning process

① Data collection and preparation - *

② Feature selection - As well as identifications of features that are useful for the learner, it is also necessary that the features can be collected without significant expense or time, and that they are robust to noise and other corruption of data that may exist in the collection process.

③ Algorithm choice

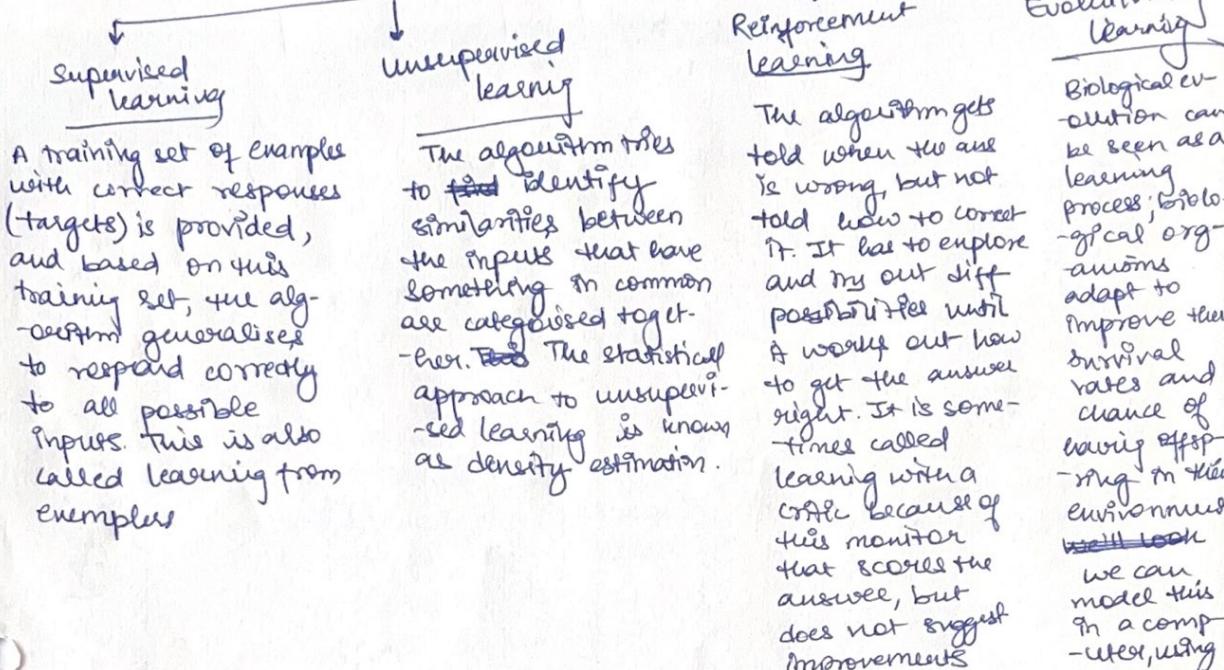
④ Parameter and model selection - for many of the algorithms there are parameters that have to be set manually, or that require experimentation to identify appropriate values.

⑤ Training - Given the dataset, algorithm and parameters, training should simply the use of computational resources to build a model of the data in order to predict outputs on new data

⑥ Evaluation - Before a system can be deployed it needs to be tested and evaluated for accuracy on data that it was not trained on.

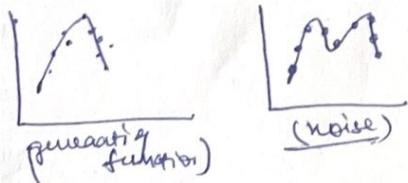
* machine learning algorithms need significant amounts of data, ~~and~~ preferably without too much noise, but with increased ~~dataset~~ size comes ^{increased} ~~excessive~~ computational costs

* Types of machine learning system



* Testing ML Algorithms

* overfitting - If we are training for too long, then we will overfit the data, which means that we have learnt about the noise and inaccuracies in the data as well as the actual function. Therefore, the model that we learn will be too much complicated, and won't be able to generalise.



We want to stop the learning process before the algorithm overfits, which means that we need to know how well it is generalising at each timestamp. We can't use the training data for this because we will not be able to detect overfitting, but we can't use the testing data either because we are saving it for the final tests. So we need a third set of data to use for this purpose, which is called the validation set because we are using it to validate the learning so far. This is also known as cross-validation in statistics. It is a part of model selection.

* Training, Testing and validation sets

3 sets of data:

- ① training set to actually train the algorithm
- ② ~~test~~ validation set to keep track of how well it is doing as it learns
- ③ ~~test~~ set to produce final results.

The semi-supervised learning attempts to deal with this need too large amounts of ~~data~~ labelled data.

Each algorithm is going to need some reasonable amount of data to learn from (precise needs vary, the more the data the algorithm sees, the more likely it is to have seen examples of each possible type of input, although it also increases the computational time). However the same argument can be used to argue that validation and test sets should also be reasonably large.

Generally, the exact portion of training: testing: validation is up to you, but it is typical to use 50:25:25 when sufficient data is available, and 60:20:20 otherwise.

In order to split the data, it is advised that randomly reordering the data is done first, or by assign each datapoint randomly to one of the sets

If you are short of ~~data~~ training data, so that if you have a separate validation set there is a worry that algo will not be sufficiently trained; then it is possible to perform leave-one-out, ~~multiple~~ multi-fold cross-validation.

Dataset is randomly partitioned into k subsets and one subset is used as validation while algo trains on the ~~others~~ all the others. A different is then left out and a new model is trained on that subset, repeating the process on all subsets. Finally the model with lowest Validation error is tested and used. We've traded off

data for computation time. In the most extreme case of this there is leave-one-out cross-validation, where the algo is validated on just one piece of data, training on all the rest.

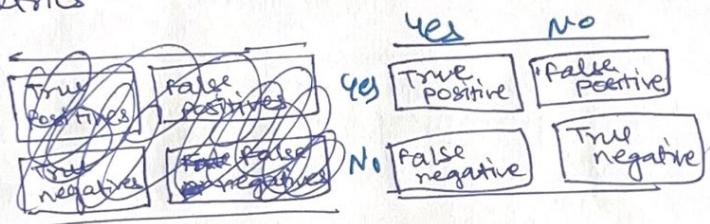
* confusion matrix - It is a method suitable for classification problems, to work out whether or not the result is good make a square matrix that contains all possible classes in both horizontal and vertical directions and list the classes along the top of the table as predicted outputs and then down to the left-hand side as the targets

The element (i,j) tells how many input parameters were put into class i in target j but class j in algorithm. Anything on the leading diagonal is the correct answer.

outputs			
	C_1	C_2	C_3
C_1	5	1	4
C_2	0	4	2
C_3	0	0	4

Here accuracy is the sum of elements on the leading diagonal divided by the sum of all elements in the matrix.

* Accuracy metric



The entries on the leading are correct.
Accuracy is sum of no. of TP and TN divided by sum of no. of examples total
(Based on binary classifications)

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

$$\text{sensitivity} = \frac{\# \text{TP}}{\# \text{TP} + \# \text{FN}}$$

$$\text{specificity} = \frac{\# \text{TN}}{\# \text{TN} + \# \text{FP}}$$

$$\text{Precision} = \frac{\# \text{TP}}{\# \text{TP} + \# \text{FP}}$$

$$\text{Recall} = \frac{\# \text{TP}}{\# \text{TP} + \# \text{FN}}$$

Two complement pairs to interpret the performance of a classifier, namely sensitivity and specificity, and precision and recall.

$$F_1 = 2 \frac{\text{Precision} \times \text{recall}}{\text{Precision} + \text{Recall}}$$

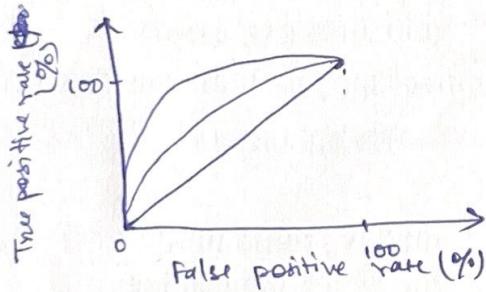
$$\text{or } F_1 = \frac{\# \text{TP}}{\# \text{TP} + (\# \text{FN} + \# \text{FP})/2}$$

* Receiver operator characteristic (ROC) curve

We can compare classifiers - either with the same classifier with different learning parameters, or completely different classifiers. In this case, ROC is useful.

This is a plot of the percentage of true positive on y-axis against false positive ~~on~~ on the x axis. A perfect ~~classified~~ classifier would be at point (0,1), while the anti-classifier that got everything wrong would be at (1,0).

Any classifier that sits on the diagonal line from (0,0) to (1,1) ~~ted~~ behaves exactly at the chance (assuming that the positive and negative classes are equally common)



anything above the line is better than the chance, and further from the ~~ted~~ line the better.

In order to compare classifiers, or choices of parameters settings for the same classifier, you could just compute the point that is furthest from the chance line on the diagonal. However it is normal to compute the Area under curve (AUC) instead.

The key to getting a curve rather than a point on ROC curve is to use cross-validation. If you use 10-fold cross-validation, then you have 10 diff classes, with 10 diff test sets and you also have the ground truth tables. The truth table can be used to produce a ranked list of diff cross-validation-trained results which can be used to specify a curve through the 10 points on ROC that correspond to results of this classifier. ~~for~~ ~~producing~~

By producing ROC for each classifier it is possible to compare the results

* Unbalanced datasets

Balanced dataset → when there are same no. of positive & negative examples in dataset, ~~otherwise~~
otherwise unbalanced dataset

$$\text{Balanced accuracy} = \frac{\text{sensitivity} + \text{specificity}}{2}$$

~~Matthews~~

Matthew's correlation coefficient, MCC

$$= \frac{\#TP \times \#TN - \#FP \times \#FN}{\sqrt{(\#TP + \#FP)(\#TP + \#FN)(\#TN + \#FP)(\#TN + \#FN)}}$$

If any of the brackets in denominator are 0, then the whole denominator is set to 1. This provides balanced accuracy computation.

In case of more than 2 classes, the calculations become too complex, however, it is possible that you use one class as positive and everything else as negative, and repeat this for each of different classes.

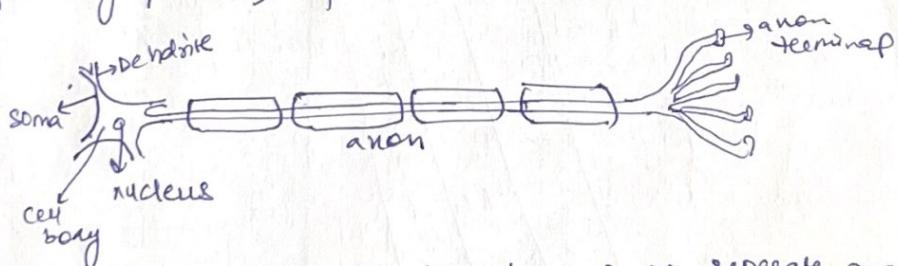
* Measurement precision

If we feed in a set of similar outputs in algorithm, then we should expect similar outputs. This measure of variability of the algorithm is also known as precision and it tells how repeatable the predictions that algo. makes are. Just because an algo is precise, does not mean it is accurate; it may be precisely wrong. One measure of how well the algo's predictions match reality is known as trueness, and it can be defined as the avg. distance b/w correct output and the prediction.

Unit 2

* The Brain and the neuron

Neurons are the processing units of the brain. There are 100 billion (10^{11}) neurons and all of different types depending upon their task. However, their general operation is similar in all cases: to transmit chemicals within the fluid of the brain raise or lower the electrical potential inside the body of the neuron. If this membrane potential reaches a certain threshold, the neuron spike or fires and a pulse of fixed strength and duration is sent down the axon. The axons divide (arborise) into connections to many other neurons, connecting to each of these neurons in a synapse. Each neuron is typically connected to thousands of other neurons. Estimated 100 billion = 10^{14} synapses in the brain. After firing, the neuron must wait for some time to receive its energy (the refractory period) before it can fire again.



Each neuron can be viewed as a separate processor performing a simple computation: deciding whether or not to fire.

How does learning occur? The principal concept is plasticity: modifying the strength of synaptic connections between neurons and creating new connections.

* Neural networks onwards in notes.

* PCA steps:

- ① Normalize data
- ② covariance matrix $\text{cov}(x,y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N-1} = \frac{\sum xy_j}{N-1}$
- ③ calculate eigen values and eigen vectors
- ④ order the eigen values from largest to smallest
- ⑤ form feature vector
- ⑥ form principle components

Unit 3

* Decision tree:

Supervised learning algo used to build classification and regression models.

A decision tree is a tree where each-

① Node - a feature (attribute)

② Branch - a decision (rule)

③ Leaf - an outcome (categorical or continuous)

* ID3 (Iterative Dichotomised 3):

$$\text{Entropy, } H(S) = \sum_{c \in C} -p(c) \log_2 p(c)$$

$$\text{Information gain, } IG(A, S) = H(S) - \sum_{t \in T} p(t) H(t)$$

Steps: ① calculate entropy for dataset

② for each attribute/feature

- calculate entropy for all its categorical values

- calculate IG for each feature

③ Find feature with maximum information

④ Repeat until we get the desired tree

For continuous attributes, we need to create a new boolean value that is true when humidity $\leq c$ and false otherwise.

Calculate best threshold of c :

① sort examples
② value of c that maximizes IG must lie at changing boundaries

③ Find IG for all such values

④ Select best attribute.

$$H(\text{value}) = ?$$

$$H(\text{feature}) = \sum_{\text{all values of feature}} p(\text{value}) H(\text{value})$$

$$IG = H(S) - H(\text{feature})$$

Characteristics of IG:

① uses greedy approach; can get stuck in local optima

② can overfit to training data

③ not always smallest tree

④ handle to use on continuous data.

* CART (Classification & Regression trees):

- is a decision tree methodology

- umbrella term for following types of DT:

① Classification trees
② Regression trees

In classification:

- Gini index is a metric for classification tasks in CART. It stores sum of squared probabilities of each class

$$Gini = 1 - \sum_{i=1}^n p_i^2$$

steps: ① calculate gini index for each feature

② select feature with lowest gini index

$$G(\text{feature}) = \sum_{\text{all values of feature}} p(\text{value}) H(\text{value})$$

Unit 4

clustering

* measuring sim. & dissimil.

① Distances:

→ Manhattan distance,

$$d_{ij} = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots$$

→ Euclidean:

$$d_{ij} = \sqrt{(x_{i1} - x_{j1})^2 + \dots}$$

② Binary vectors:

m_{01} = the number of attribute with
p 0 and q was 1

m_{11}

m_{10}

m_{00}

Simple matching coefficient

$$SMC = \frac{m_{11} + m_{00}}{m_{11} + m_{00} + m_{10} + m_{01}}$$

~~Jaccard~~

Jaccard coefficient,

$$J = \frac{m_{11}}{m_{01} + m_{10} + m_{11}}$$

③ * Types of clustering

partitioning approach

- ① K-means
- ② K-medoids
- ③ CLARANS

↓
Hierarchical approach

- ① Dianal (divisive)
- ② Agnes (agglomerative)
- ③ BIRCH
- ④ ROCK
- ⑤ CAMELEON

density-based approach

- ① DBSCAN
- ② OPTICS
- ③ Pendule

* Centroid, Radius & diameter of cluster

Centroid (middle of circle), $C_m = \frac{\sum_{i=1}^N f_i p_i}{N}$

Radius (sq. root of avg dist from any point to centroid), $R_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (f_i p_j - f_m)^2}{N(N-1)}}$

Diameter, $D_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (f_i p_j - f_m)^2}{N(N-1)}}$

③ Cosine similarity!

If d_1 and d_2 are two document vectors, then

$$\text{cos}(d_1, d_2) = \frac{d_1 \cdot d_2}{\|d_1\| \|d_2\|}$$

$$\|d_1\| = \sqrt{\sum x_{i1}^2}$$

⊗

* k-means clustering: Partition clustering approach

Pseudo code:

1. select k points as initial centroids
2. repeat
3. form k clusters by assigning all data points to the closest centroid
4. recompute the centroid of each cluster
5. until centroids don't change

- ① → Initial centroids chosen at random
 - closeness is measured by Euclidean dist, cosine similarity, correlation, etc.
 - often stopping condition changed to "until relatively few points change"
 - complexity is $O(n * k * I * d)$
 - ↑ num of points
 - ↓ no. of clusters
 - ↓ no. of iterations
 - no. of attributes
- ② → Relatively higher strength. ~~than PAM and CLARA~~ $O(I * k * n)$ than PAM and CLARA
- ③ → often terminates at local optima. The global optima may be found using ~~deterministic~~ annealing and genetic algorithms

Disadv

- ① Need to specify k in advance
- ② unable to handle noisy data and outliers
- ③ not suitable for clusters with non-convex shapes
- ④ has problems when clusters are of differing sizes, densities and non-globular shape

Evaluating k-means:

$$SSE = \sum_{i=1}^k \sum_{x \in C_i} \text{dist}^2(m_i, x)$$

sum of squared errors

one easy way to reduce SSE is to increase k

* hierarchical:

- used distance metric as clustering criteria
- method does not need no. of clusters as input, but needs a termination condition.
- represented in tree like structure, called dendrogram

Type → Agglomerative - known as bottom up approach
consider it as bringing things together

Dissolveive - (top down approach) we take a large cluster and divide it into two, three, four or more clusters

~~Agglomerative~~

Agglomerative

Stopping criteria:

- ① Pick several clusters (k) upfront
the next merge will create a cluster with low cohesion.
- ② Stop when
- ③ a. Diameter of cluster
b. Radius

Divisive:

- * more efficient than agglomerative clustering especially when there is no need to generate a complete hierarchy all the way down to individual level.
- * can be considered as global approach, since it contains the complete info. before splitting data.

Issues in Divisive:

- ① Splitting criteria - Ward's k-means square error is used here (SSE)
The greater reduction obtained in difference in SSE should reflect goodness of split.
Gini Index ~~for~~ for handling nominal data
- ② Splitting method - bisection k-means approach
- ③ choosing cluster to split - by checking the square errors of clusters and splitting the one with largest value.

K-means vs Hierarchical

- ① Hierarchical cannot handle big data, k-means can (because order of k-means is linear $O(n)$, while that of hierarchical is $O(n^2)$)

- ② In k-means, since we start with random choice of centers, the results produced might differ. While results in hierarchical are reproducible
- ③ K-means works well with hyper spherical shaped clusters

- ④ K-means require prior knowledge of k , but not in hierarchical

* Partitioning method

This clustering method classifies the information into multiple groups based on characteristic and similarity. In partitioning methods, when database(D) that contains multiple (N) objects then the partitioning method constructs user-specified (k) partitions of data in which each partition represents a cluster and a particular region.

- Eg: ① k-means
 ② PAM (k-medoids)
 ③ CLARA (Clustering Large Applications)

PAM → find representative objects, called medoids, in clusters.
 : built in 1965 by Kaufman and Rousseeuw

Pseudo code:

- ① Select k representative objects arbitrarily
- ② For each pair of non-selected object h and selected object i,
 calculate the swapping cost T_{ih}
- ③ For each pair of i and h,
 - if $T_{ih} < 0$, i is replaced by h
 - then assign each non-selected object to the most similar representative object
- ④ Repeat step ②-3 until there is no change

* Total swapping cost, $T_{\text{sum}} = \sum_j C_{jh}$

$$C_{jh} = d(j, h) - d(j, i)$$

* PAM is more robust than k-means in the presence of noise and outliers because medoids are less influenced by outliers and works efficiently for small datasets but does not scale well for large datasets.

$$\Theta(k(n-k)^2)$$

CLARA: It draws multiple samples of data, applies PAM on each sample and gives best clustering as output.

Strength: deals with large datasets

Weakness: efficiency depends on sample size

A good clustering on sample does not ensure good clustering of the whole dataset if the sample is biased.

* Distribution-based methods: It is a clustering model in which we will fit the data on the probability that how it may belong to the same distribution. The grouping of data may be normal or gaussian. Gaussian is more prominent where we have a fixed no. of distributions and all the incoming data is fitted into it such that the distribution of data may get minimized.

This model works well on synthetic data and diversely sized clusters.

Problems: ① If the constraints are not used to limit the model's complexity.
 ② produces clusters that assume concisely defined mathematical model underlying the data, rather strong assumption for data distribution.

- * Density based
 - * based on density (local cluster criteria), such as density connected points.
- * major features
 - ① Discover clusters of arbitrary shape
 - ② Handle noise
 - ③ One scan
 - ④ Need density parameters and as termination condition
- * Two parameters:
 - ① Eps : max. radius of neighbourhood
 - ② minpts : min. no. of points in an ~~Eps~~ neighborhood of real point
- * $N_{\text{Eps}}(P) : \{q \text{ belongs to } D \mid \text{dist}(P, q) \leq \text{Eps}\}$
- * Directly density-reachable: A pt is directly density-reachable from point q w.r.t $\text{Eps}, \text{minpts}$ if
 - ① P belongs to $N_{\text{Eps}}(q)$
 - ② core point condition:
 $|N_{\text{Eps}}(q)| \geq \text{minpts}$

- * ~~Density~~ density-reachable : A point p is density-reachable from a point q w.r.t $\text{Eps}, \text{minpts}$ if there is a chain of points $P_1, P_2, \dots, P_n, P_1 = q, P_n = p$ such that P_{i+1} is directly density-reachable from P_i
- * Density connected : A point p is density-connected to a point q w.r.t $\text{Eps}, \text{minpts}$ if there is a point o such that p and q are density-reachable from o w.r.t $\text{Eps}, \text{minpts}$.

- * ~~DBSCAN~~ DBSCAN (Density based spatial clustering of Applications of Noise):
 - Relies on a density-based notion of cluster: A cluster is defined as a maximal set of density-connected points
 - Discover clusters of arbitrary shape in spatial databases with noise.
- Algorithm:
 - ① Arbitrarily select a point P
 - ② Retrieve all points density-reachable from P w.r.t Eps and minpts
 - ③ If P is a core point, a cluster is formed
 - ④ If P is a border pt., no points are density-reachable from P , and DBSCAN visits next point of database.
 - ⑤ Continue this process until all of the points have been processed

* Fuzzy clustering

It is a type of clustering algorithm in machine learning that allows a data point to belong to more than one cluster with different degrees of membership. (between 0 and 1)

Applications: ① Img segmentation, pattern recognition, risk assessment, medical diagnosis.

Advantages ① Flexibility
② Robustness
③ Interoperability

Disadv: ① Complexity
② Model selection

Algorithm: ① Initialise the data points into the desired no. of clusters randomly

- ② Find out the centroid ~~as $\sum_{k=1}^m \sum_{i=1}^n x_{ik}$~~
- ③ find distance of each pt from centroid
- ④ update membership values
- ⑤ Repeat steps 3-4 until constant values are obtained for membership values or the difference is less than a tolerance value
- ⑥ Defuzzify the obtained membership values.

centroid formula, $v_{ij} = (\sum_{k=1}^n v_{ik}^m x_k) / \sum_{k=1}^n v_{ik}^m$

~~as $\sum_{k=1}^n v_{ik}^m$~~
m is fuzziness parameter

Updating membership values,

$$v_{ik} = \left[\sum_{j=1}^m \left(d_{ki}^2 / d_{kj}^2 \right)^{\frac{1}{m-1}} \right]^{-1}$$

* Evaluation of unsupervised learning models

① 2 classes of techniques to validate results → internal validation
→ external validation

* Internal validation revolves around two metrics:

① cohesion within each cluster

$$\text{cohesion} (C_k) = \sum_{x_i, y_i \in C_k} \text{similarity}(x_i, y_i)$$

② separation b/w different clusters

$$\text{separation}(C_j, C_k) = \sum_{x_i \in C_j, y_i \in C_k} \text{distance}(x_i, y_i)$$

A set of clusters with high cohesion & high separation are considered ~~very~~ good.
Some measures that combine both of the above into a single measure:

- ① Silhouette coefficient
- ② Calinski - Harabasz coeff.
- ③ Dunn index
- ④ Xie - Beni score
- ⑤ Hartigan index.

* External validation - This can be carried out if true cluster labels are available. The idea is to measure the statistical similarity to the true cluster set.

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

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

$$F_1 \text{ factor} = \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

TP = No. of pairs of records which are in same cluster for both P and S

FP = No. of pairs which are in same cluster in S but not in P

* Twin sample validation: How to validate the results of our unsup. learn. model in absence of true cluster labels. This step takes it as a given that we have already performed clustering on our training data and now we want to validate the results.

- * Steps:
- ① Create twin sample of ~~the~~ training data
 - ② Performing unsup. learn. on twin sample
 - ③ Importing results for twin sample from training set
 - ④ Calculating similarity b/w two sets of results