# AML
## LAB ASSIGNMENT - 2

PROBLEM STATEMENT - To study feature engg and implement the dimensionality reduction techniques ( PCA & TSNE)

### OBJECTIVES -

① To understand feature engg & learn feature selection techniques
② To understand the concept of dimensionality reduction.

### THEORY -

→ Feature Engg
- Process of transforming raw data into features that better represent the underlying problem to the predictive models, resulting in improved model accuracy on unseen data.
Some feature engg methods are -
① Binning
② log transform
③ Scaling
④ Feature selection.

→ Feature selection techniques
- Process where you automatically or manually select those features which contribute to your prediction value or o/p in which you are interested in. This method can deal w/
① Multicollinearity

② High date dimension
③ High training time

feature selection methods
① Forward selection
② Backward elimination

⟶ Dimensionality reduction (PCA, LDA t-SNE)

- Filter method
① ANOVA - Analysis of variance is a method used to analyze the diff among the group means in a sample

② Chi square - This is defined as where $O_i$ is the observation & $E_i$ is the expected value - This value of chi can be used to derive the p-value that gives us the probability of independence. If p value is high ($> 0.05$), we can say that the attr is not statistically significant to the target var.

③ Pearson's correlation - This is defined as The value is the measure of strength of linear association b/w 2 variables, where $r = 1$ means a perfect +ve correlation and the value $r = -1$ means a perfect -ve correlation.

- Wrapper method
① RFE - It is a wrapper type feature selection
   algo. This means that a different ML algo
   is given & used in core of the method
   & wrapped by RFE and used to help select
   features

- Intrinsic method
① DT - A decision tree is a method which can be
   traversed based on the attr values and
   can give an intrinsic value at the leaf
   nodes.

- PCA
* Algo that uses the eigen values derived from
  the correlation metrix in order to reduce
  the dimension of the dataset. The reduced
  features are representative of dataset but
  doesn't hold any meaning on its own.

- t-SNE
* t - distributed stochastic gradient descent
  neighborhood embedding is ML algo that
  employees stochastic neighbor embedding
  to reduce the no of attr by projecting
  them on low dimension space.

CONCLUSION — feature Engg was studied. Implemented the two dimensionality reduction techniques PCA & t-SNE.

## FAQs

**Q(1)** What are the various dimensionality reduction techniques?

**Ans)** The various dimensionality reduction techniques are —

① Ratio of missing values
② Low variance in the column values
③ High correlation b/w 2 columns
④ Principal component analysis (PCA)
⑤ Candidates and split columns in random forest
⑥ Backward feature elimination
⑦ Forward feature construction
⑧ Linear discriminant analysis (LDA)
⑨ Neural autoencoder
⑩ t-distributed stochastic neighborhood embedding (t-SNE)

**Q2)** Define feature engg?

**Ans)** Process of using domain knowledge to extract features from raw data via data mining techniques. These features can be used to improve performance of ML algos.

Q7) Difference b/w t-SNE & PCA ?

| PCA | t-SNE |
|---|---|
| ① Linear dimensionality reduction technique | ① Non linear dimensionality reduction technique |
| ② Tries to preserve the global structure of data | ② Tries to preserve local structure (cluster) of data |
| ③ Deterministic algo | ③ Non deterministic or randomised algo |
| ④ Works by rotating the vectors for preserving variance | ④ Works by minimising the distance b/w the point in a gaussian |