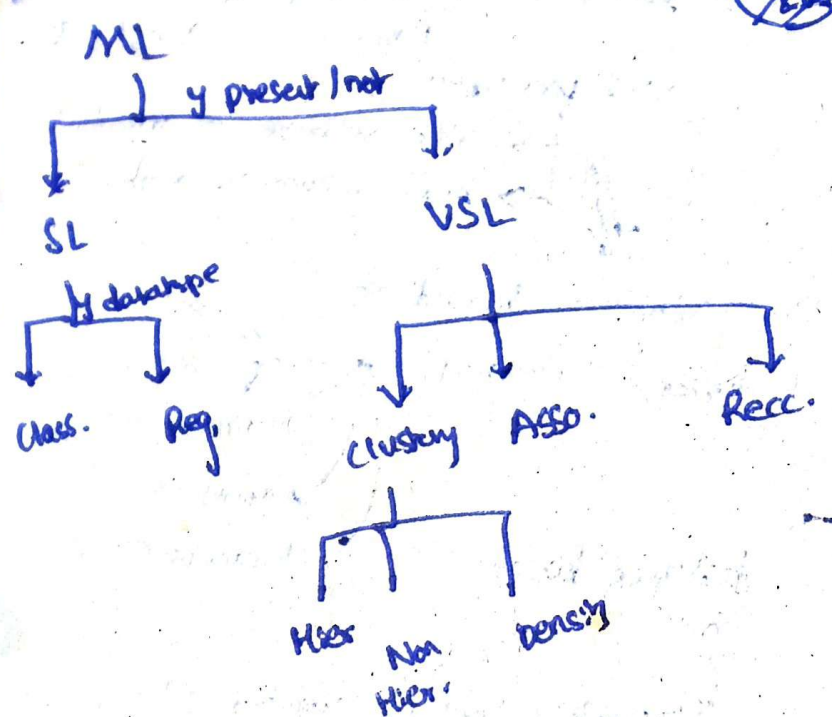


SLI 2



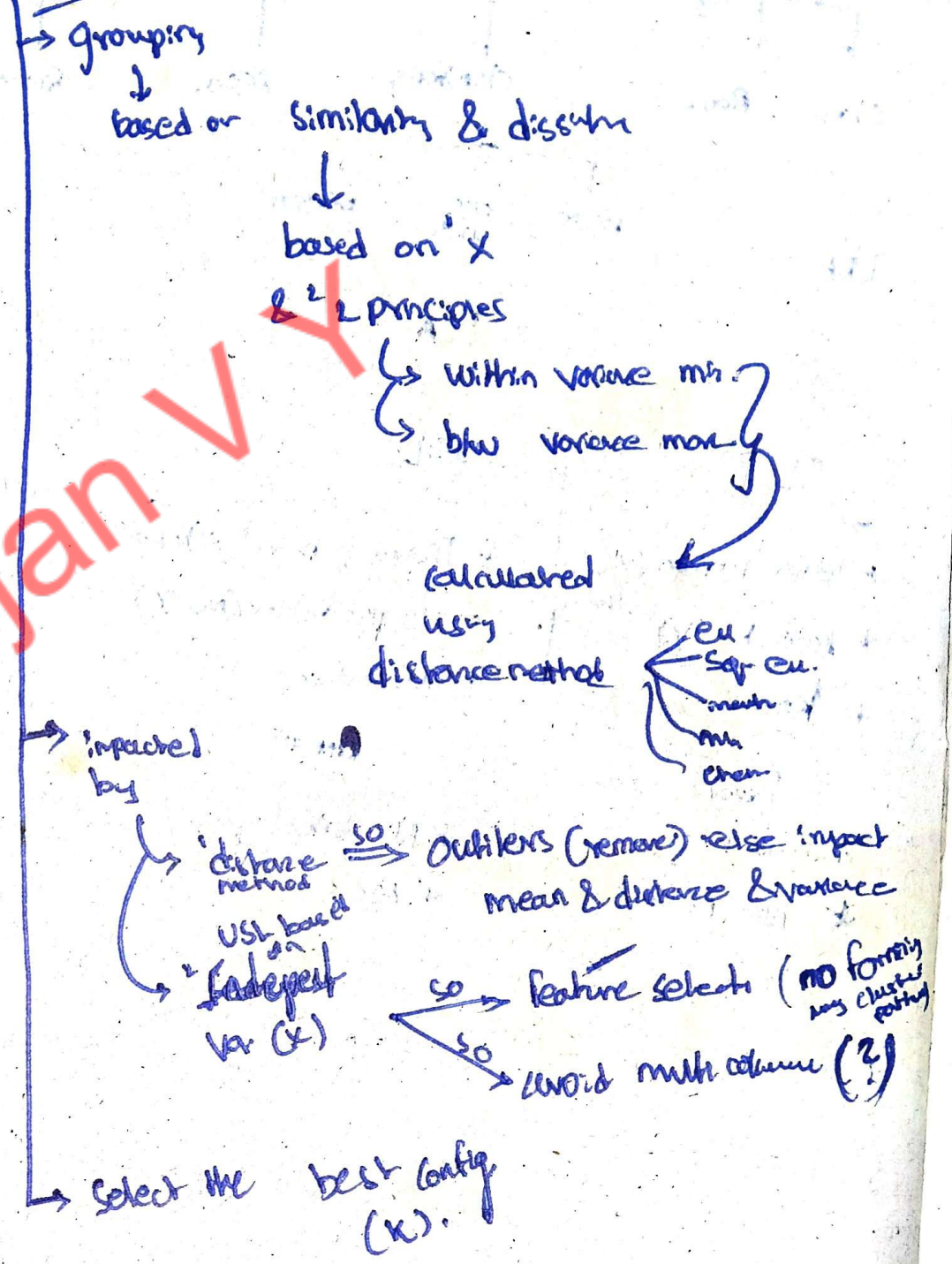
LDA
→ SL feature/
dimension
reduction

PCA → VSL feature/dimension
reduction

- * ✓ Pattern (x/feature)
- * ✓ Predict (y)
- * ✓ Metric (y)
- * ✓ Pattern (x/feature)
- * x predict (no y)
- * x metric (no y)

known (label) vs unknown (feature) stage

Clusters what?



↳ How? / Algo-2.

↳ k means clusters.

no. of clusters
↓
specifies mean of cluster

1. Choose k (data no.)
2. Choose distance method (eu.)
3. align datapoints to centroid
4. align centroid to centroid
5. align datapoints to centroid
6. do 4,5 until convergence.

→ Job of centroid is to align to the centre of the datapoints in the group

datapoint
Shift blue group
↓
centroid changes
(means blue groups)
↓
centroid changes.

Stop / convergence
↳ no datapoint
Shift blue groups

clustering framework of varying.

1. Define Business problem
(FD, DC, IC, etc.)

↓
2. Feature select

↓
3. choose distance method

↓
4. choose cluster method

↓
5. Build cluster

↓
6. profile / naming cluster
by datapoint (not ago.)

↓
7. effectiveness

Silhouette score

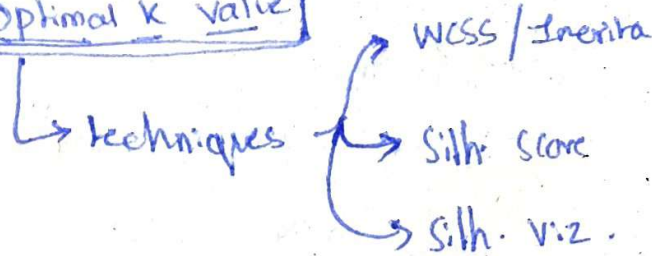
no. of datapoint
... ..

Silh root value

algorithm

complete
average
center

Optimal K value



WSS

what? within cluster sum of squares

$$WSS_{k=N} = WSS_{cluster} + \dots$$

--- + WSS cluster N

variance / distance
b/w datapoint &
Centroid

no. of features

$$\sum_{i=1}^n (x_i - \bar{x}_i)^2$$

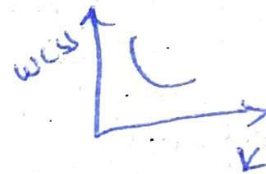
How to use it to find wing it?

↑ k → ↓ WSS

as datapoints
in clusters
align more closer
to the centroid.
↳ so reduce
variance / distance

How to use it?

elbow plot



what: axis: k y-axis: WSS
elbow shape (ideally)

How:

drastic decrease in WSS
& then flatten out

limitation: subjectively not
drastic decrease
& flatten out

so, elbow plot is mainly used to eliminate the
unnecessary k value: e. to that k.

Silh. score

→ based on 2 principles

within var. distance b/w datapoint & centroid is min. (a)
b/w var. distance b/w datapoint & centroid is max (b).

→ what?

$$\text{Silh. score} = \frac{\sum_{i=1}^{\text{no. of datapoint}} \text{Silh. coeff}}{\text{No. of datapoints}}$$

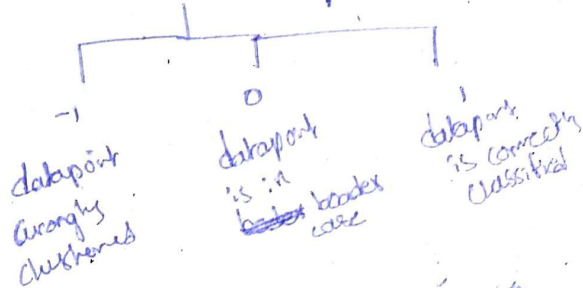
→ what is Silh. coeff?

$$\text{Silh. coeff.} = \frac{(b-a)}{\max(b,a)}$$

(is w.r.t datapoint)

b → b/w var. distance b/w datapoint & centroid mean.
a → within cluster distance b/w the point & other points mean.

→ what Silh. coeff. signifies?



$$\text{Silh. coeff. of datapoint} = \frac{(b-a)}{\max(b,a)} = \frac{(a-a)}{\max(a,a)} = \frac{0}{a} = 0$$

$$\text{Silh. coeff. of datapoint} = \frac{(b-a)}{\max(b,a)} = \frac{(0-a)}{\max(0,a)} = \frac{-a}{a} = -1$$

$$\text{Silh. coeff. of datapoint} = \frac{(b-a)}{\max(b,a)} = \frac{0}{b-a} = 0 \text{ if } b=a$$

Silh. viz

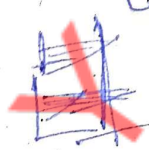
→ 3 criteria (for choosing optimal k value)

1. no. outliers in any cluster (i.e. no. of Silh. coeff. in any cluster)
2. Silh. coeff. of most of datapoints is greater than the Silh. score
3. Highest Silh. score

→ How to Silh. viz?

1) For each cluster

→ Silh. coeff. arranged in descending order



& then plotted & joined

2) Finally a vertical line to mark Silh. score.

X-axis: Silh. coeff. value.
Y-axis: cluster name.

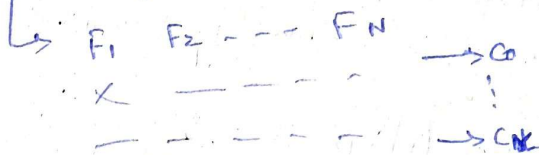
Optimal k value

1. Pair plot (only > 2 features)
2. WCSS / Inertia / Elbow plot
3. Silh. Score
4. Silh. viz

Profiling clusters

→ Non technical; business based; domain knowledge needed.

→ More related to the descriptive Analysis of the clusters.



→ then used in targeting the product ⇒ value added to the business.

→ What Silh. score signifies?

- ≈ 1 good clustered as majority of Silh. coeff. is ≈ 1.
- ≈ 0 bad clustered
- ≈ -1 poorly clustered

So, Silh. score higher i.e. closer to 1 → that k value is preferred as clustered good.

Hierarchical clustering

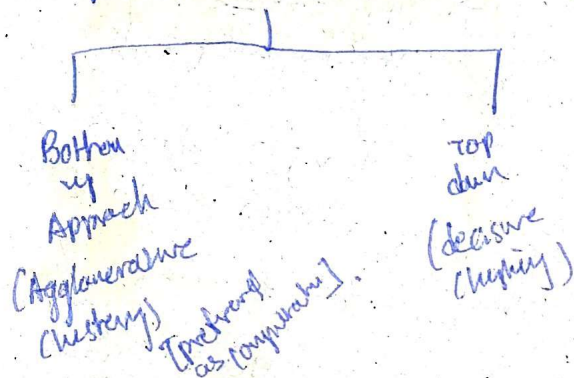
→ what is hierarchical? (based for Hs. clustering)

↑ hierarchical level → ↓ no. of persons

⇓

1) in the hierarchical clustering

→ Types of Hierarchical clustering



↑ Iteration → ↓ no. of datapoints cluster

↑ Iteration → ↓ no. of clusters (datapoints)

→ How? / Algo.

1. Each datapoint as cluster
2. pair wise distance b/w clusters
3. Merge the two nearest clusters.
4. step 2-3 until become one larger cluster.

→ as clusters (not datapoint) use linkage method for distance calculation:
[No euclidean — which for b/w datapoints]

1. Single linkage → min dist
2. Complete linkage → max dist
3. Centroid linkage/mean → b/w centroid dist
4. Avg linkage → Avg of dist.

Linkage method distance

- Ward's linkage

Out of all possible pair combination of the cluster, the cluster formation with min. variance is merged.

→ called min-variance method.

distance b/w all points in C_1 & C_2 to centroid of C_1 & C_2 .

→ i.e. centroid using all datapoints from C_1 & C_2

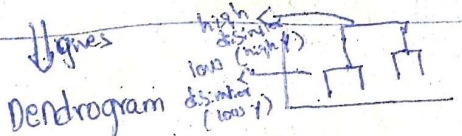
→ why? (preference) / when? to use

- compact & sphere structure
- not effect by outlier
- so, no chain effect
- Objective is min. variance.

Linkage

→ gives 3 things

1. Index of cluster(s) / datapoints to be merged.
2. Cophenetic distance / i.e. distance @ which they are merged.
3. no. of datapoints in that formed merged cluster



→ is visual representation of how cluster(s) & datapoints merged.

X-axis: Index of cluster / datapoint
Y-axis: Cophenetic distance / dissimilarity

→ optimal no. of clusters is the range in Y-axis that for which large range exist → cut @ any distance (Y value) of that range (distance/Y)

Significant distance

→ γ is dissimilarity & it is stable for long distance. → Natural group by avoid heterogeneity

Cluster formation method

→ one of disadvantage of algorithmic clustering / hierarchical clustering
→ it stop when one big cluster is formed.

→ 2 types

- distance based Cophenetic distance i.e. Y-axis in dendrogram
- max-cluster

By this two methods we form the desired no. of clusters formation.

How to decide optimal no. of clusters? (in hierarchical clustering)

Cophenetic correlation

→ is Pearson correlation b/w distance matrix & cophenetic matrix.

→ Highest value is better cluster formation. → why?

→ because euclidean distance b/w datapoints (closeness in level of datapoints) given by distance matrix need to be similar to Cophenetic distance (i.e. linkage distance) b/w datapoints (@ level of clusters)

More similar means clusters are well formed else not clustered the datapoints into right cluster.

→ Key components

- 1. Distance matrix (b/w datapoints - euclidean distance)
- 2. Cophenetic matrix (b/w cluster(s) & datapoint(s) - Cophenetic distance / linkage distance) @ the no. of clusters.

PCA

when to drop a column

1. Multi collinearity
2. NO correlation (x-y)
3. Redundancy
4. P > 0.05 - not significant
5. Regularization (lasso)
6. Domain knowledge - drop column
7. Std = 0.
8. Image → processing speed as ↑ dimensions

But info loss due to drop column.

PCA \rightarrow reduce those impacts w/o dropping column.

PCA \rightarrow based on
& 1. Information
& 2. Variation/Variance

✓ Variation	\rightarrow	✓ Info
✗ Variation	\rightarrow	✗ Info

\hookrightarrow key rule.

Curse of dimensionality

\hookrightarrow Guideline

$n \rightarrow$ no. of feature
at least $n^2 \rightarrow$ no. of rows / datapoints

to avoid curse of dimensionality

\hookrightarrow for same no. of data points

$\rightarrow \uparrow$ no. of feature
 \downarrow leads to

\uparrow sparsity

\downarrow leads to

overfit of model
(as less ^{no. of} datapoints
but more no. of feature)

\hookrightarrow is memorizes
; not find pattern
to generalize)

\hookrightarrow overcome by

PCA : (w/o dropping columns)

What is PCA?

1. Principal component analysis
↳ is used dimension reduction technique w/o dropping column.

2. Principal component(s) are new axis formed.

3. PC is linear combination of all features.

$$PC_1 = a_1x_1 + \dots + a_nx_n$$

$$PC_2 = b_1x_1 + \dots + b_nx_n$$

PC_n =

4. No. of PCs = No. of features / columns.

5. PCs are orthogonal to each other

6. PCs are vectors

↳ constraints

1. Magnitude

↳ given by

eigen value

↳ also called
Explained variance

↳ Captures the
variances / variation

(i.e. information)

2. direction

↳ given by

eigen vector

↳ gives the
direction of
their PCs
wrt original
axes

(i.e. Range
(max - min)
for
feature
column)

7. PC₁ always captures ^{most} maximum variance.

Order of variance / information / variation capturing is

$$PC_1 > PC_2 > \dots$$

8. Based on 'information', 2 branches.

Why PCA?

1. dimension reduction w/o dropping columns.
2. Reduce Noise.
3. Reduce multicollinearity

When to use

1. Only for Numerical columns / variables.
2. only if have multicollinearity

How?

Working in graphical

1. plot the data
2. Scale data → origin
Shift to mean
3. Rotate the axis till
capture max. variance
(PC₁)

4. draw other PCs
orthogonal to the
PC₁.

5. drop the PCs that
capturing least variance
/ info / variation using
the needful methods.

Algorithmic steps

1. Standardize the data
($\mu=0$, $std=1$)
2. Find covariance matrix.
3. Find eigen vector &
eigen value. $AV = \lambda V$
4. Sort the eigen values
in desc.
5. Select the Eigen vectors
based eigen value i.e.
PCs capturing maximum
Amtr of variance / variation /
information
6. Get the PCs by
dot product of eigen
vectors & standardized
data.

$$A \rightarrow \text{Covariance matrix (is square matrix)} = \frac{S^T \times S}{(n-1)}$$

(vector)

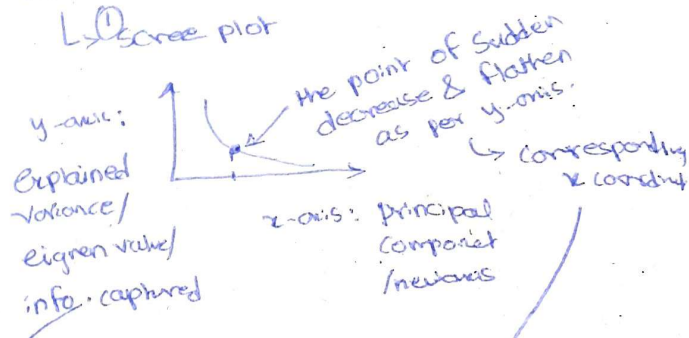
$V \rightarrow$ eigen vector (vector)

$\lambda \rightarrow$ eigen value (scalar)

How? → selecting ~~best~~ optimal no. of PCs

Three methods

↳ ① Scree plot



is the required optimal no. of PCs.

Limitation: Subjective

↳ ② Kaiser method

↳ PCs with eigenvalue (λ) > 1 are selected.

↳ based on empirical experiment that if $\lambda > 1$ then we get atleast 60% Explained variance / variation in total variance / variability

↳ Limitation: if most of $\lambda < 1$ then it is not working as optimal.

↳ ③ percentage rule (most widely used)

↳ Select PCs when arranged in descending order based on λ such that cumulative sum $\geq 80\%$.

Preference:

Percentage rule $>$ Kaiser rule $>$ Scree plot

Result / outcome

* dimension reduction w/o dropping columns / features

↳ we drop PCs based on Variance captured (λ) wrt total Variance (≤ 1)

↳ we not drop features.

* So, new dataset we work with is columns as PCs, not original features / axis.

* For test data (if only linked to Supervised learning; as no testing in UDL),

↳ we just transform based on fit on train data.

↳ so that no data leakage b/w train & test.

Two ways

