# DBSCAN Clustering

$o \rightarrow$ outlier

$\rightarrow$ border points

$\rightarrow$ close d point

$\downarrow$

Density is very high

$\epsilon = \underline{3\ mm}$ (user)

Density low

C1

1) hierarchy

2) distance

3) density — closeness → radius

$\downarrow$

draw one circle ← $\epsilon$

$\epsilon$ = radius

Cluster points

within

cluster

1) Cluster point      →   $\in$ draw ⊙

                     ↳ point within ⊙

      ↓

   ( density )

2) Border point  → point on ⊙ boundary

                     ↳ draw ⊙ with radius $\in$

      ↓

  ( extension )          ↳ Repeat this process to cover max data points

3) Noise or outliers

DBSCAN → density based algorithm

DBScan → Density Based Spatial Clustering
of Applications with Noise

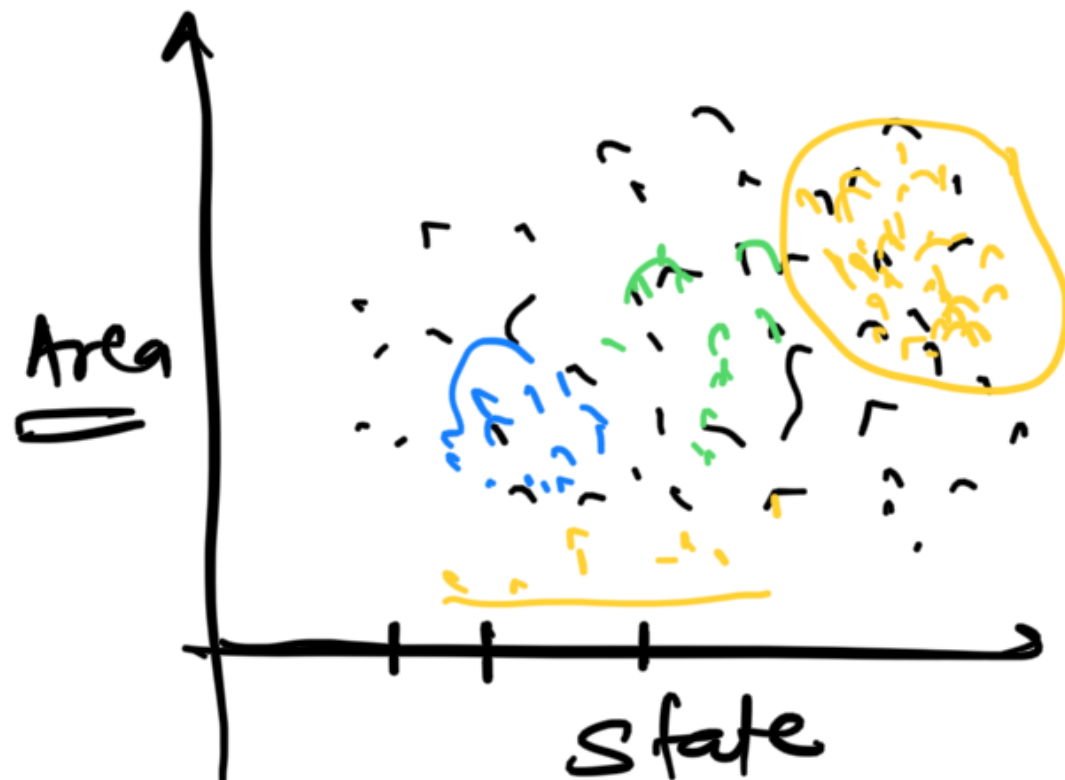App → Statewise density of languages

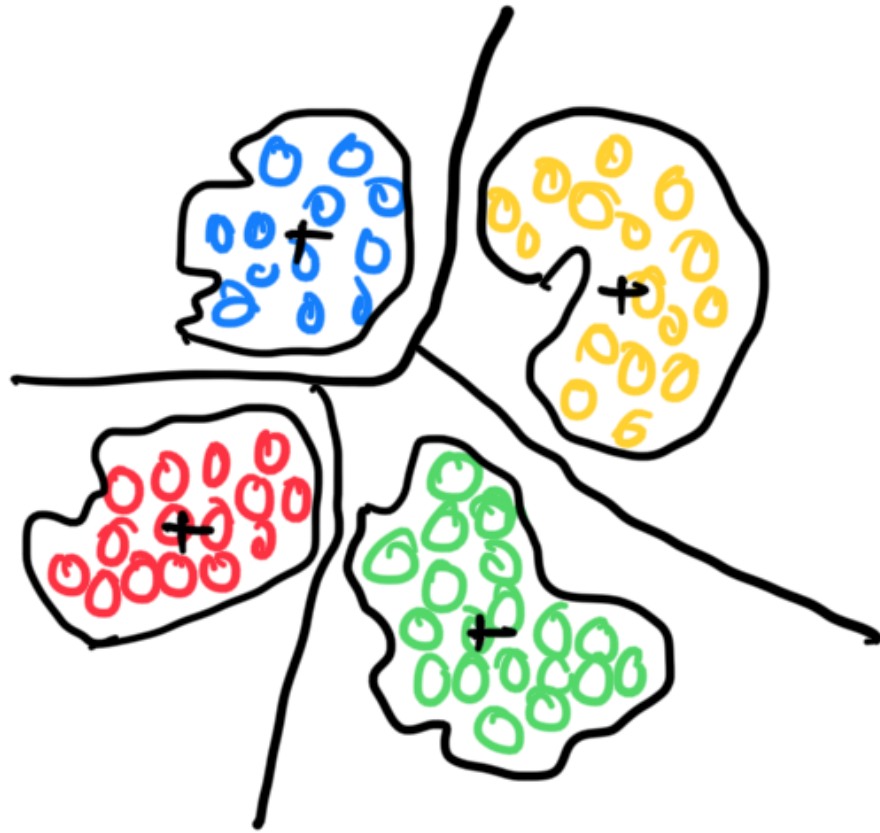→ DBSCAN

Language → 5 → cluster

→ Marathi

→ Hindi

Bengali

Tamil

Gujrati

$\boxed{\text{Density}}$ = number of points within the radius of Circle

$(\epsilon)$

$\underline{\text{Basic Idea}} \rightarrow$ Cluster the dense region

DBSCAN $\rightarrow$ Idea of density
                              ↓
                          low/high



$+$

$\boxed{\text{K Means}}$

↳ Avg — mean

↳ Medoid → median ✓
                    ↓
                centroid

$\boxed{\text{K-medoid Algorithm}}$

Values

Min - Max

d.f. describe()

min, max, r, v

Range → 0 -100 ⇒ mean

0, 200, 789, 450, 13, ⇒ medoid.

→ Mode → Mode

K - mode

Mean = 5, 3, 4, 21

$$\frac{5+3+4+21}{5}$$

=

Medoid = 100, 99, 5, 7

# DBCCAN - Concepts

1. core points - data points → within 🙂

        └→ cluster points

        └→ Radius (ε) ⇒ ε values should be
                                    same

2. Border points - Neighborhood points

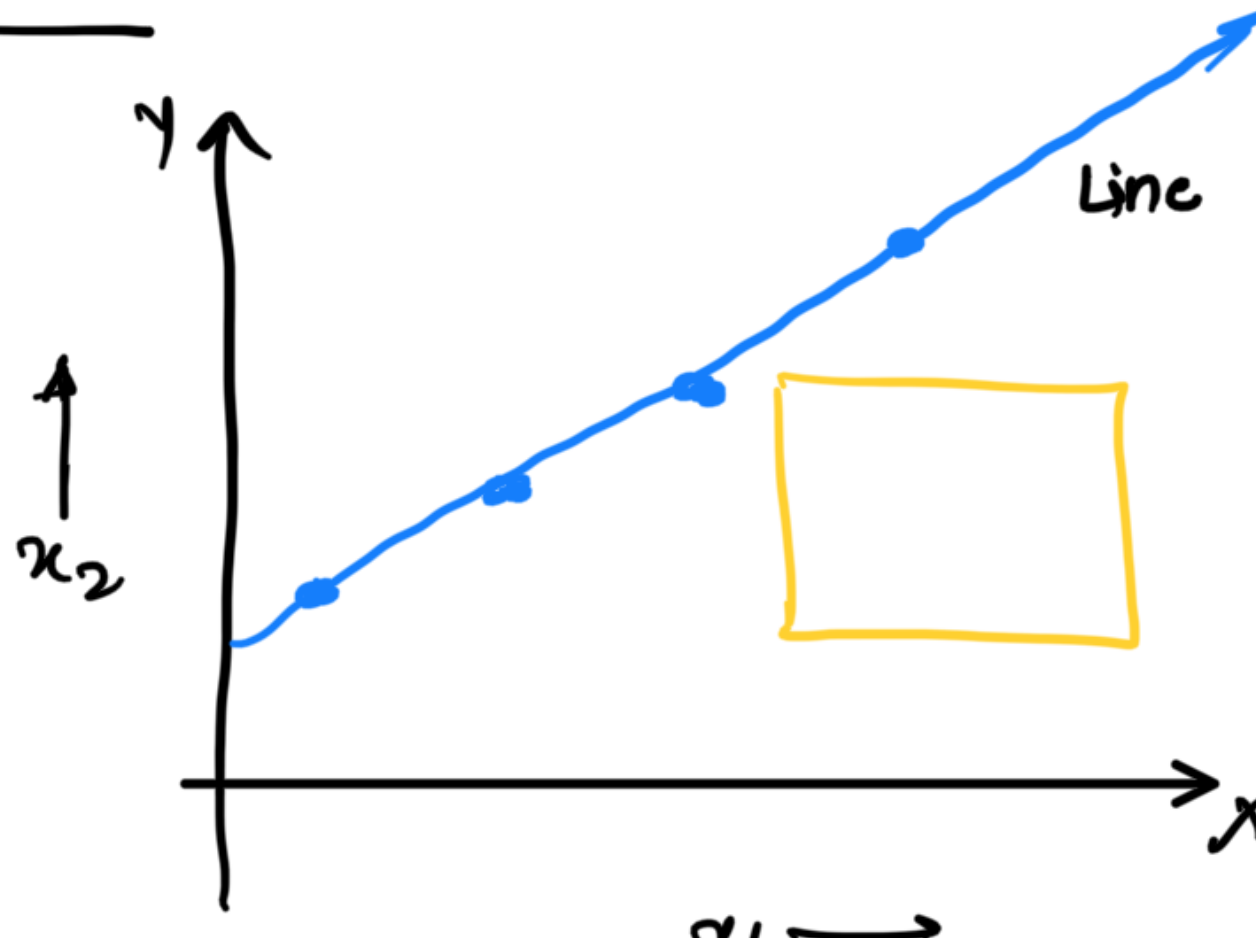3.   Noise — Outliers → Discards $\{-1\}$

DBSCA
K-Means

---

## Dimension Reduction

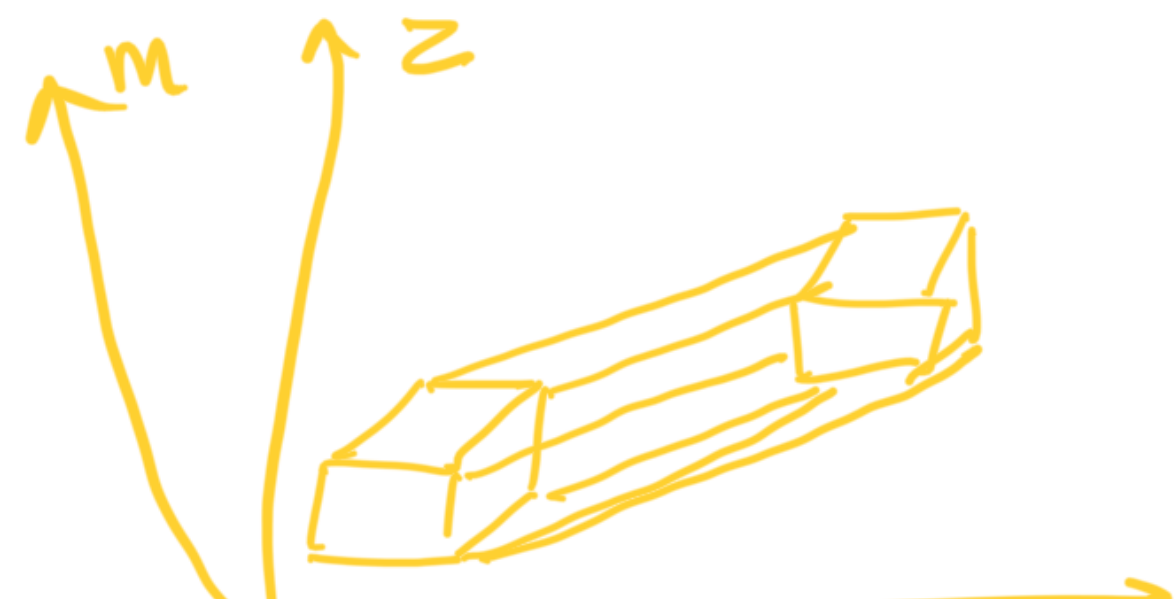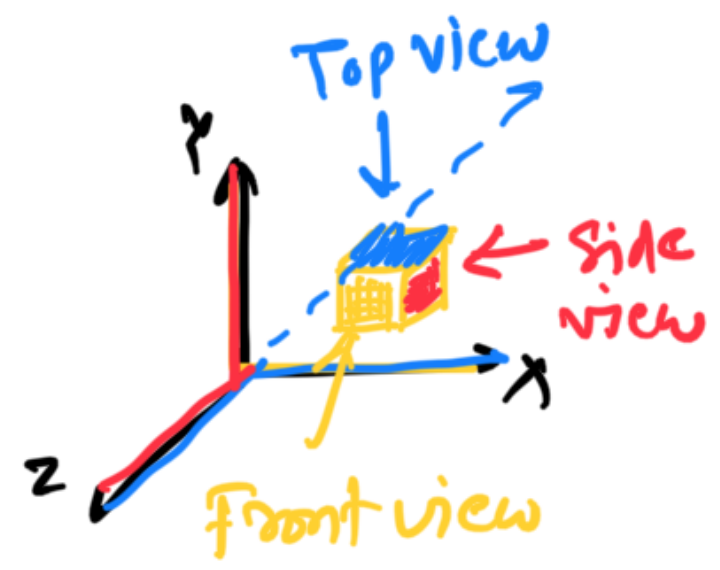| $x_1$ | $x_2$ |
|-------|-------|
| 1 | 2 |
| 2 | 4 |
| 3 | 6 |
| 4 | 8 |

Line

2D

plane

3D

Top view

Side view

Front view

y

x

z

m

z

10, 100, 1000 → visualization high dimension
is very difficult
↓
Complexity
↑
Imagination

# Dimension Reduction

Ex:
- Documents → multiple page
  (words)     ↳ multiple words
              ↳ multiple letter

Alphabets → 2̶6̶ X 300 X 15 X 30 X 6 X 12

= ⟨ x ? ⟩
      ↑        words
   Huge number  ↓
                features

| the | for | sunset | — | — | — | — | — | — | n |
|-----|-----|--------|---|---|---|---|---|---|---|
|     |     |        |   |   |   |   |   |   |   |
|     |     |        |   |   |   |   |   |   |   |

Lakhs

P1 →
P2 →
P3 →
.
.
Pn →

Difficult task

Image ⇒ ?

High dimensional data

daily
man
words

P1 P2 P3 ... P15

Ex. 2   Image

450
↓
x  300

Bank ⇒ SBI →

↓

105 → centers in Country

25 → States

10 → City

→ In each center

↳ 10,000

$$105 \times 25 \times 10 \times 10000 \Rightarrow \text{Customers} \quad ?$$

↓

Huges of data

Eg. Images → Traffic Control ⇒ Cameras → 10L

↓
24×7
↓
1 Lac

$$10L \times 24 \times 7 \times 1 \, Lac$$

= Image dataset

= 1

1 Camera → 1 day → 1 Lac

7 days # 7 lac

$10 L \times 7 Lac \Rightarrow 70 Lacs$



1 sec
1 sec

video

Eg. Genetics $\rightarrow$ Gen. $\rightarrow$ $G_1, G_2$



2D

3D

$S_{10}$

$g_2$

$G_1 \quad g_1 \quad g_1 \quad \cdots \quad G_{200}$ → Dimensions more

|       | $G_1$ | $g_1$ | $g_1$ | | $G_{200}$ |
|-------|-------|-------|-------|---|---------|
| $S1$  |       |       |       |   |         |
|       |       |       |       |   |         |
| $S10K$ |      |       |       |   |         |

1) meaningful visualisation

2) how to extract that meaning informations

Lung Cancer → ?

Tumor →    > 10000

↑

Stage of
Tum

2,50,000 → each cell.

$2,50,000 \times 10000$

→ $Sol^n$ ⇒ Dimension Reduction

1) Convert your high dimensional data
into lower dimensional data

    a) Visualize the data

    b) Analyze the data

Dimension Reduction $\longrightarrow$ loss of relevant inform.

minimum loss of
relev. info
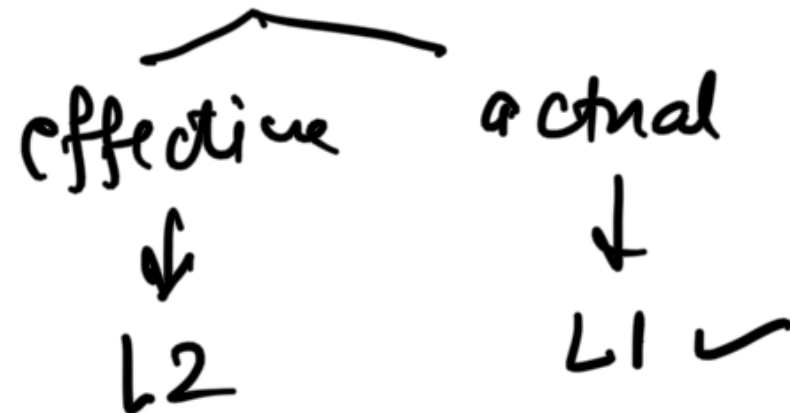
## Approaches to DR -

1) Feature Selection

    — Select subset of existing features

2) **Model Regularization** ✓

- L2 reduces → dimensionality

effective     actual
↓             ↓
L2            L1 ✓

3) **Combining of existing feature into**
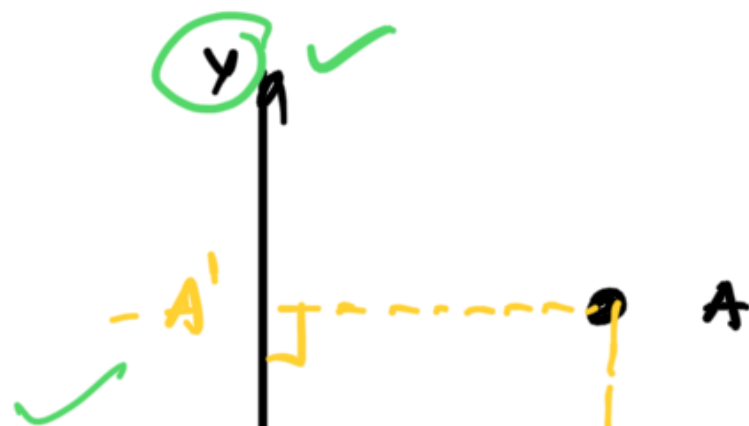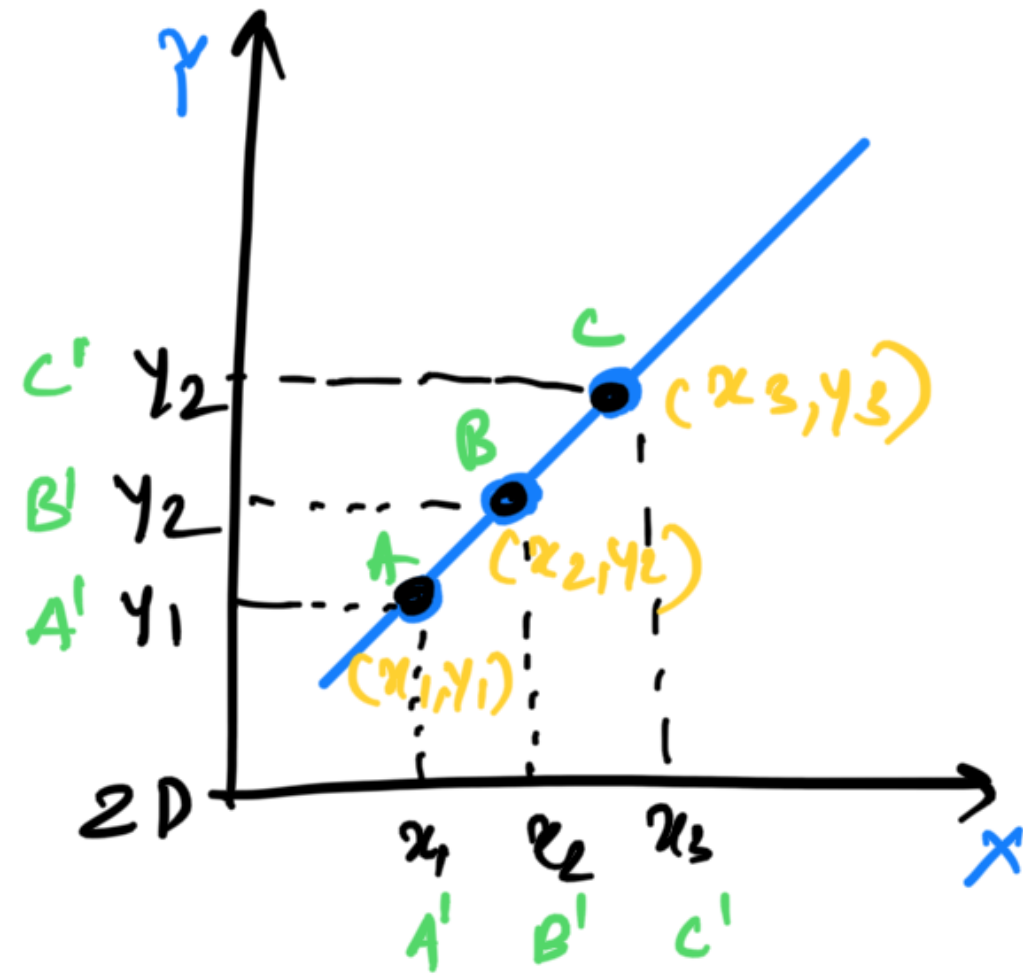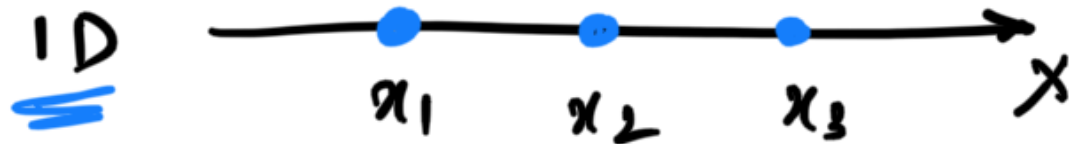**smaller no. of new features**

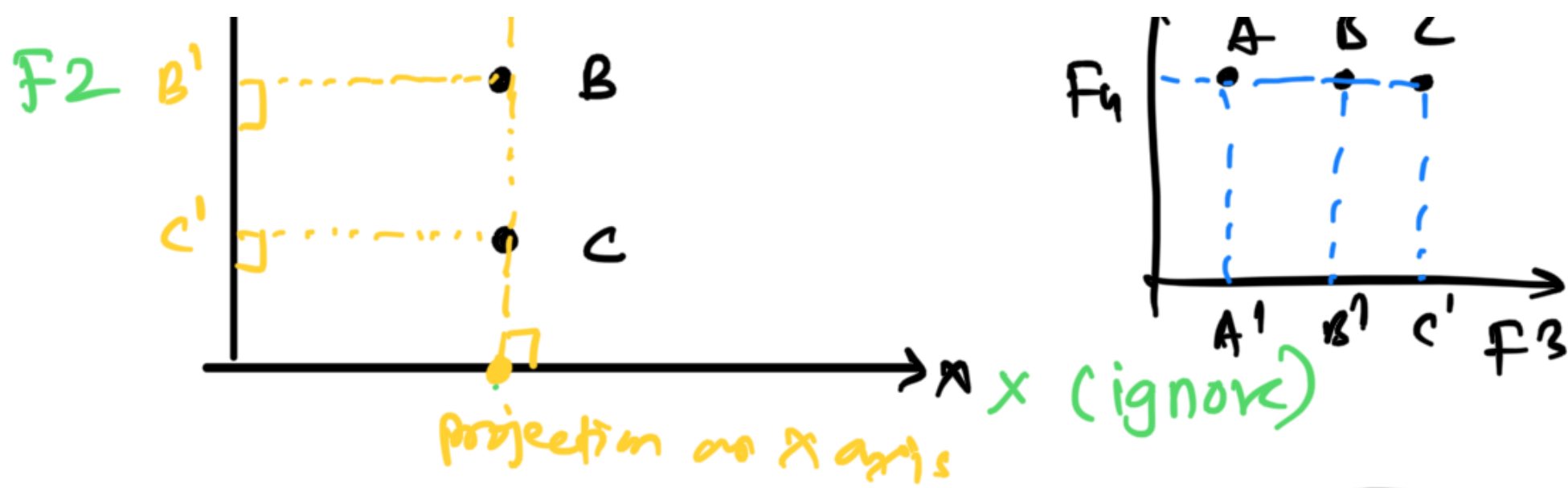— Projection ( Linear combination)

— LDA, PCA, Kernel PCA

## ② Combining of existing features into new feature

### Projection



1D

2D

$c'$ $y_2$ — $C$ $(x_3, y_3)$

$B'$ $y_2$ — $B$ $(x_2, y_2)$

$A'$ $y_1$ — $A$ $(x_1, y_1)$

$x_1$ $x_2$ $x_3$

$A'$ $B'$ $C'$

Ⓨ ✓

$A'$ ┅ $A$

Linear Discriminant Analy

LDA

F2   B'                          B

     C'                          C

                              → X  x (ignore)
Projection on X axis

( F1 (Ignore) )  X      ( F3 )
     F4                 ( F4 )

F4   A  B  C

        A'  B'  C'  F3

# PCA – Principal Component Analysis

— linear projection — 90°

PC2                    x3

PC1        x1

Y                   All pts
                    we are
                    getting

        Y    x

                    x1 & x2

projection          Rotate   PC1 → PC1

x2    A'

$x_3 \Rightarrow x_1'$

$x_1 \oplus x_2$

$x_8'$

$x_1'$

$x_1$

PC2

$y$

$x$

PC1

2 comp

$x_1 \oplus x_2 \longrightarrow x_1' \sim PC1$

$x_3 \oplus x_4 \longrightarrow x_2' \sim PC2$

$$x_1 + x_2 = \frac{x_1 + x_2}{\text{Day1 Math}}$$

Day1  Math

$x_1'$ Vs $x_2'$

$x_5'$

→ Polynomial → PCA
Regr

ICA