

→ Dimensionality Reduction →

① PCA [Principal Component Analysis]

Consider

→

<u>weight</u>	<u>height</u>	<u>DBP</u>	<u>SBP</u>	<u>Health?</u>

↓
Diastolic B.P.

↓
Systolic B.P.

Say Health of person depends on 4 features

- weight ✓
- Height ✓
- DBP ✓
- SBP ✓

To represent this data ⇒ We need 4 Dimension

For more features → We need more Dimension.

* Visualizing data in More than 3 dimension is difficult.

* Computation on 4 features is also complex

Possible Solution →

<u>Height</u>	<u>weight</u>	<u>DBP</u>	<u>SBP</u>	<u>Health?</u>

↑ ↑ ↑ ↑

[Body Mass Index] BMI

BP [Blood Pressure]

① look for strong correlated features: Here Height & weight } are strongly correlated
Also DBP and SBP } correlated

... correlation of Height & weight ⇒ BMI

→ We can combine the effective correlation of Height & Weight \Rightarrow BMI
 " " " " " " " " DBP and SBP \Rightarrow BP.

Now Instead of 4 columns we have 2 columns.

BMI	BP

→ This is 2-D data.
 Easy to Compute
 Easy to Visualize.

We have Reduced the dimension of data from 4-D to 2-D.

Note \Rightarrow Consider 2-columns:-
 let \rightarrow X_1 X_2

DBP	SBP
✓78	✓126
80	128
81	127
82	130
84	130
86	132

$$y \text{ (BP)} = \alpha_1 \text{ (DBP)} X_1 + \alpha_2 \text{ (SBP)} X_2$$

α_1 & α_2 are weights of features
 X_1 and X_2 are features

Ex \rightarrow ① $BP = 0.8 \text{ DBP} + 0.6 \text{ SBP}$

We are giving more weights to DBP as compared to SBP.

Ex ② let $BP = \text{mean of DBP \& SBP}$

for mean

DBP	SBP	BP mean
78	126	102
80	128	104
81	127	104
82	130	106
84	130	107
86	132	109

$$BP = \frac{DBP + SBP}{2} = 0.5 \text{ DBP} + 0.5 \text{ SBP}$$

$$\alpha_1 = 0.5$$

$$\alpha_2 = 0.5$$

... and SBP.

82	150	102
84	130	107
86	132	109

Ex (3)

let BP = Sum of DBP and SBP.

$$BP = DBP + SBP \Rightarrow \alpha_1 = 1$$

$$\alpha_2 = 1$$

For Sum

DBP	SBP	Sum
78	126	204
80	128	208
81	127	208
82	130	212
84	130	214
86	132	218

↑

Now PCA \Rightarrow Principal Component Analysis ✓

\Rightarrow It is a method to find the Linear Combination that
accounts for as much variability as possible
in Combined Variable \rightarrow (Maximum Variance)

[To understand $y = \alpha_1 x_1 + \alpha_2 x_2$
we want value of α_1 and α_2 such that the variance in
the value of (y) should be maximum
B'coz move the variance \rightarrow Move is the information]

Note: \rightarrow for $y = \alpha_1 x_1 + \alpha_2 x_2$
(combined variable)

\rightarrow Here to maximize variance large value of α_1 and α_2 can be
proposed

\rightarrow We will place restriction on value of weights i.e. ($\alpha_1, \alpha_2, \dots$)
such that

$$\alpha_1^2 + \alpha_2^2 + \dots + \alpha_n^2 = 1$$

* If there are n features, we will need ' n ' weights

eg. \rightarrow we had 2 features: we need α_1 and α_2

$$\begin{aligned} \therefore BP &= \alpha_1 DBP + \alpha_2 SBP \\ \text{constraint } \alpha_1^2 + \alpha_2^2 &= 1 \end{aligned}$$

Let $\alpha_1 = 0.8$

$$(0.8)^2 + (0.6)^2 = 0.64 + 0.36 = \underline{\underline{1.00}} \quad \checkmark$$

let $\alpha_1 = 0.8$
 $\alpha_2 = 0.6$

$$(0.8)^2 + (0.6)^2 = 0.64 + 0.36 = \underline{\underline{1.0}} \quad \checkmark$$

let $\alpha_1 = 0.8$ and $\alpha_2 = 0.6$

$BP = \alpha_1 DBP + \alpha_2 SBP$

DBP	SBP	BP
78	126	138
80	128	140.8
81	127	141.0
82	130	143.6
84	130	145.2
86	132	148.0
	Mean	142.8

$$\begin{aligned} \text{variance} &= \frac{1}{n-1} \sum_{i=1}^n (BP - \text{mean BP})^2 \\ &= \frac{1}{5} \left((138 - 142.8)^2 + (140.8 - 142.8)^2 + (141.0 - 142.8)^2 + \dots \right) \\ &= \underline{\underline{12.74}} \end{aligned}$$

let for above 6 input sample.

* We can take many values of α_1 and α_2 such that $\alpha_1^2 + \alpha_2^2 = 1$ and calculate variance for each α_1 and α_2

α_1	α_2	var(y)
0.8	0.6	12.74
0.6	0.8	11.8
0.98	0.2	10.4
0.2	0.98	7.4

of the above 4 cases: we see for $\alpha_1 = 0.8$ and $\alpha_2 = 0.6$ we get largest variance = 12.74

Tazika 1 : $BP = 0.8DBP + 0.6SBP$

To find α_1 and α_2

We got our α_1 and α_2

We got our α_1 and α_2
Task 2 \rightarrow
Now \rightarrow Consider again \rightarrow

DBP	SBP
78	126
80	128
81	127
82	130
84	130
86	132
x	y

Construct Covariance Matrix =

$$\begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix}$$

$$\text{cov}(x, y) = \frac{1}{N-1} \sum (x_i - \bar{x})(y_i - \bar{y}) = \text{cov}(y, x)$$

For above sample data set

Covariance Matrix \Rightarrow

	DBP	SBP
DBP	8.17	5.97
SBP	5.97	4.97

Eigen values =

from above Covariance Matrix.

$$\begin{bmatrix} -0.8 \\ -0.6 \end{bmatrix} \checkmark$$

$$\therefore \begin{cases} BP = -0.8 \text{DBP} - 0.6 \text{SBP} \\ BP = 0.8 \text{DBP} + 0.6 \text{SBP} \end{cases} \checkmark$$

Applications \rightarrow

(1) Consider

DBP	SBP	weight	Height

BP BMI

BP	BMI

Reduced Dimension by combining features.

BP

BMI

Reduced Dimensionality
Combining features-

(2)

	DBP	SBP	Weight	Height	Pulse
Pat 1	78	126	82	5.1	90
Pat 2	80	128	74	5.4	80
	81	127	64	5.1	70
	82	130	82	6.2	90
	84	130	92	6.3	95
	86	132	85	6.10	94

Here if we have many no. of features and we need to compare the health chart/parameters of samples, it is difficult

But if we combine the above 5 col in 2 col.

C1	C2
V1	V2
V3	V4
V5	V6

Now with less no. of features comparing the health chart/parameters of samples is relatively Easy:

* PCA \rightarrow It is Dimensionality Reduction Technique.

Consider: For House Price Prediction

Size	Location	YOC	OC Y/N	Builder	Color	Garden	Swimming	Metro Distance	Price

* Here Price of House depends on above features.

* Here many no of features (Problem of Plenty).

* Issues \rightarrow Representⁿ Problem ✓
* Overfitting Problem. ✓

* We need to Reduce Dimension

* We need to Identify Important Components.

Dimensionality Reduction \rightarrow .

* Reduces the dimension of feature space

Ex if there are 100 features/col in dataset and you want to get only 10 features then with dimensionality reduction technique we can achieve this.

* It transforms dataset which is in n dimension space to n' dimension space where $n' < n$.

to n' dimension space where $n' < n$

Why Dimensionality Reduction?

↳ Normally it is argued that many features gives
More accurate result

* But after some point the model performance decreases
(overfitting) with increase in number of features.

* This is Curse of Dimensionality

* So Dimension Reduction is crucial.

* PCA enables us to identify the correlation and pattern
in the dataset so that it can be transformed into new
dataset with lower dimension without loss of important
information.