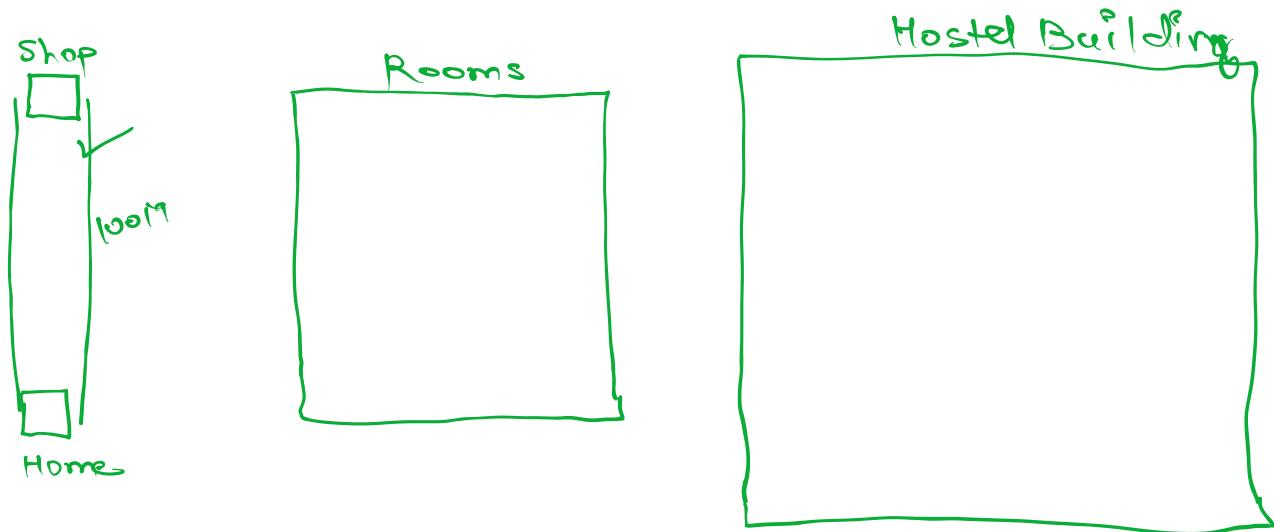


Curse of Dimensionality

No of cols = No of Dimensions



House Price Pred

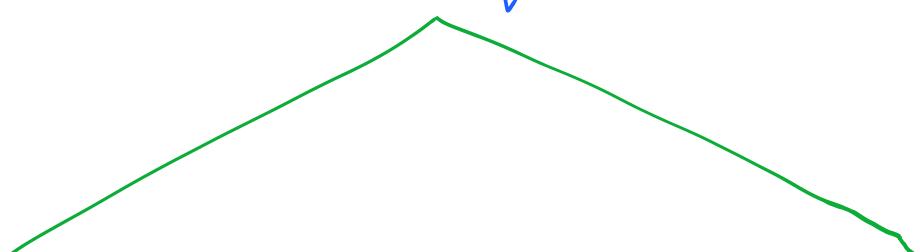
10 cols

100 cols yds

- Unnecessary cols → Overfitting
- Multicollinearity → High Dimension

Solution ?

Dimensionality reduction
Techniques (DRT)



Principal Component Analysis (PCA)

↓
Applied on all types of data

Linear / Discriminant

Analysis (LDA)

↓
only on classification data

What are the options available to us for reducing dimensions →

① Feature Elimination / feature selection

→ Hyperparameter Tuning.

→ VIF (Variance Inflation Factor) for dealing with multicollinearity.

→ RFE (Recursive Feature Elimination).

→ Tree Based Algorithms like DT, RF, XGBoost have inbuilt feature selection techniques.

② Feature Extraction

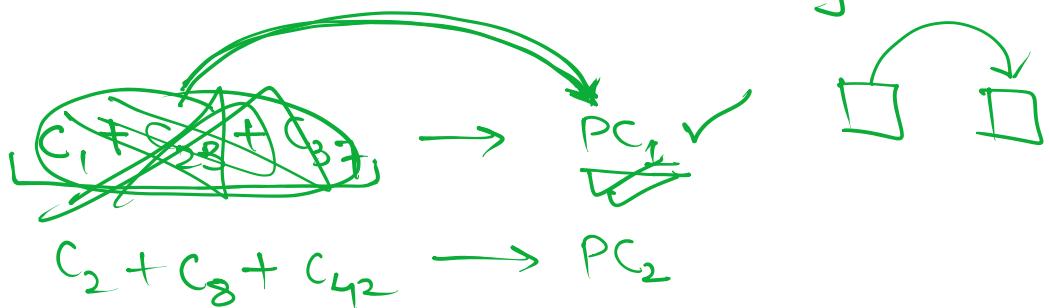
→ PCA

→ LDA

Principal Component Analysis (PCA)

$$C_1, C_2, \dots, C_{100} \xrightarrow{\text{PCA}} PC_1, PC_2, PC_3, \dots, PC_n$$

To create Principal components, PCA will combine columns linearly.



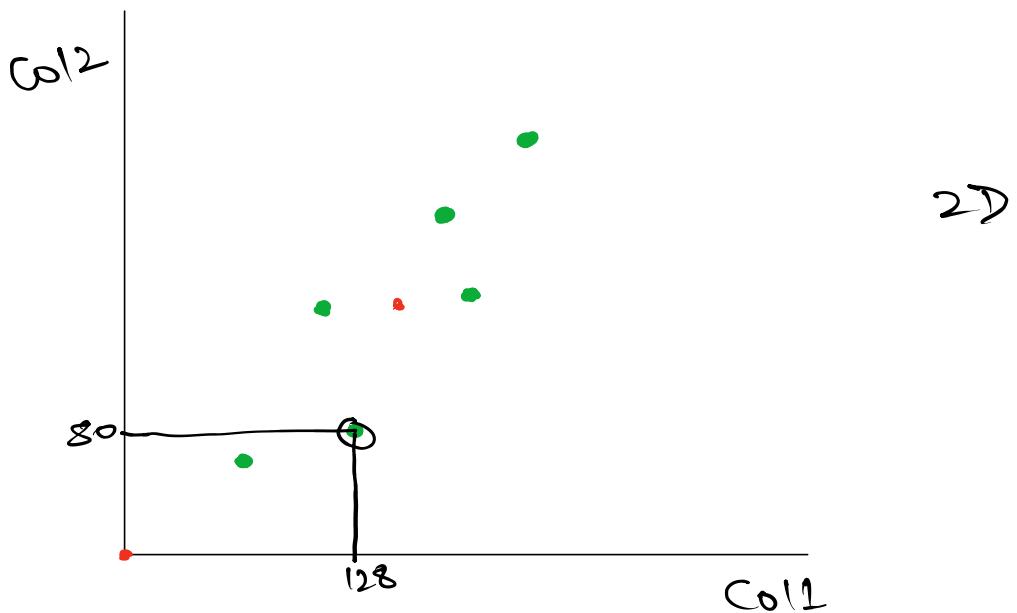
Mathematical working of PCA →

Steps in PCA →

- ① Mean centering of the data.
- ② Calculate the covariance matrix.
- ③ Calculate eigenvalues & eigenvectors of the covariance matrix.
- ④ Arrange the eigenvectors in descending order of eigenvalues.
- ⑤ Calculate the principal components.

Example calculation

Col 1	Col 2
126	78
128	80
128	82
130	82
130	84
132	86
<u>129</u>	<u>82</u>



①

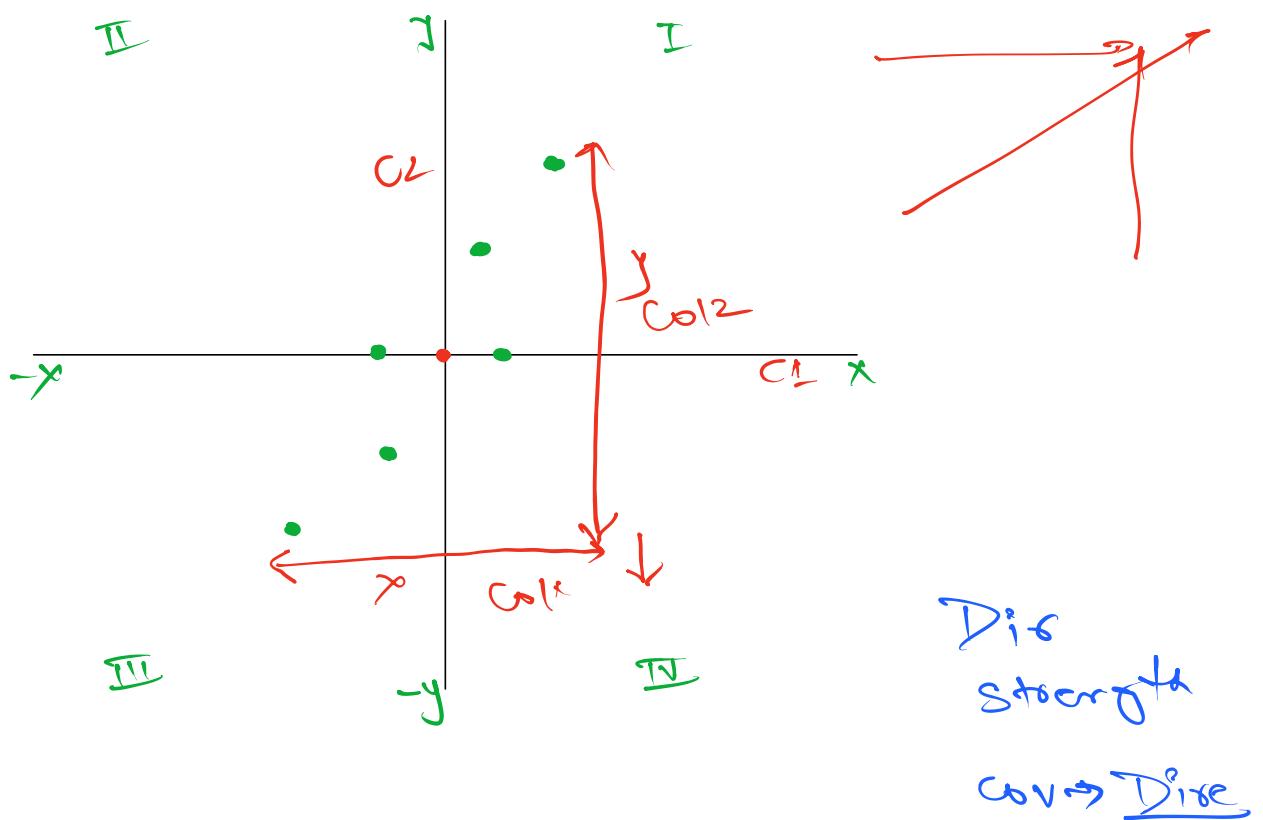
$$x - \text{mean}(x)$$

Data - Mean

Col 1	Col 2
$126 - 129$	$78 - 82$
$128 - 129$	$80 - 82$
$128 - 129$	$82 - 82$

$$\begin{array}{ccc|cc} 130 & -129 & 82 & -82 \\ 130 & -129 & 84 & -82 \\ 132 & -129 & 86 & -82 \end{array}$$

var(x)		
Col1	Col2	
$x_1 - 3$	-4	y_1
$x_2 - 1$	-2	y_2
$x_3 - 1$	0	y_3
$x_4 1$	0	y_4
$x_5 1$	2	y_5
$x_6 3$	4	y_6
<u>(0)</u> , <u>0</u>		



② Calculate the covariance matrix.

	Col 1	Col 2
Col 1	Var(c1)	Cov(c1, c2)
Col 2	Cov(c2, c1)	Var(c2)

$$\text{Cov}(c_1, c_2) = \text{Cov}(c_2, c_1)$$

Direction in which the variance of the data is more, the amount of information present in direction is more.

$$\text{Var}(c_1) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{(-3-0)^2 + (-1-0)^2 + \dots}{5} = 4.4$$

$$\text{Cov}(c_1, c_2) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) =$$

	Col 1	Col 2
Col 1	4.4	5.6
Col 2	5.6	8.0

$$\text{Cov. Matrix} = \begin{bmatrix} 4.4 & 5.6 \end{bmatrix}$$