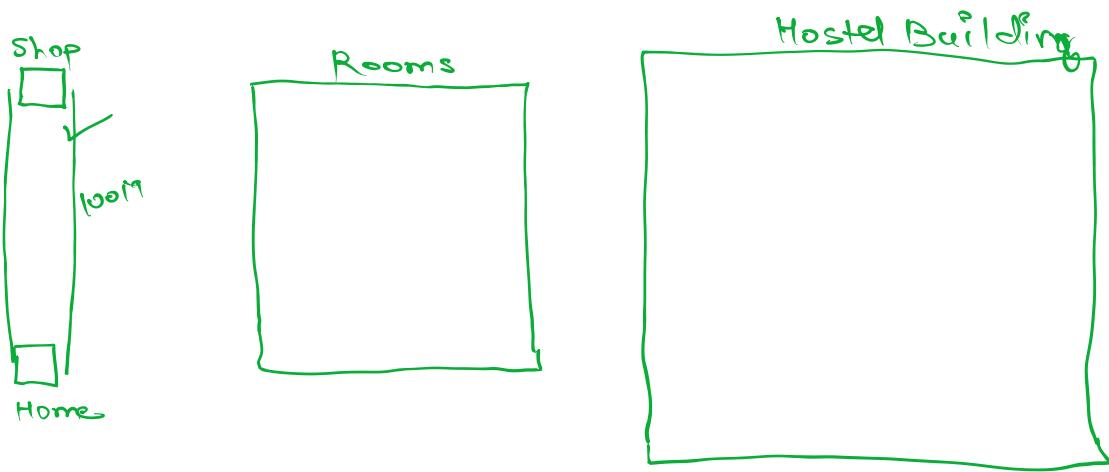


Curse of Dimensionality

No of cols = No of Dimensions



House Price Problem

10 cols

100 cols n cols

- Unnecessary cols → Overfitting
- Multicollinearity → High Dimension

Solution?

Dimensionality reduction Techniques (DRT)

Principal
Component
Analysis (PCA)

Linear
Discriminant
Analysis
(LDA)

Applied on all
types of data

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

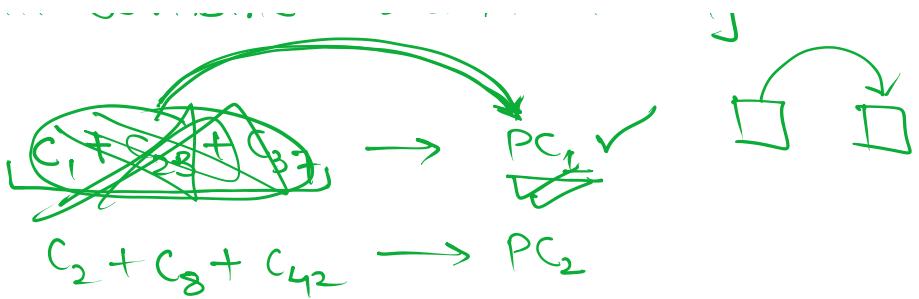
→ PCA

→ 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.

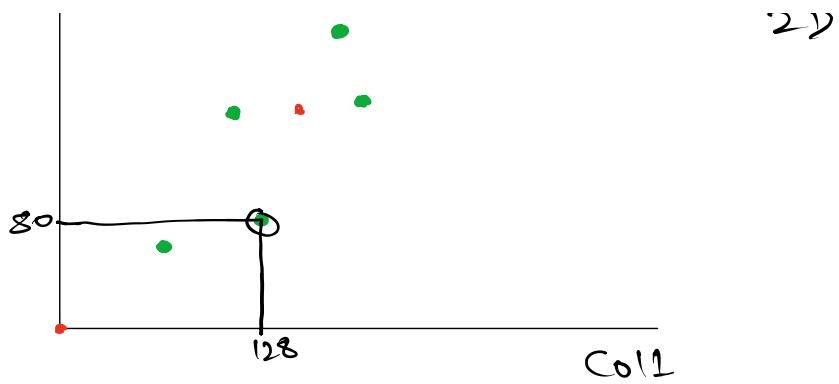
Example calculation

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

Col 2

•

- -



①

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

Data - Mean

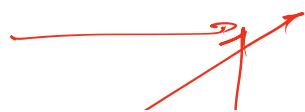
Col1	Col2
126 - 129	78 - 82
128 - 129	80 - 82
128 - 129	82 - 82
130 - 129	82 - 82
130 - 129	84 - 82
132 - 129	86 - 82

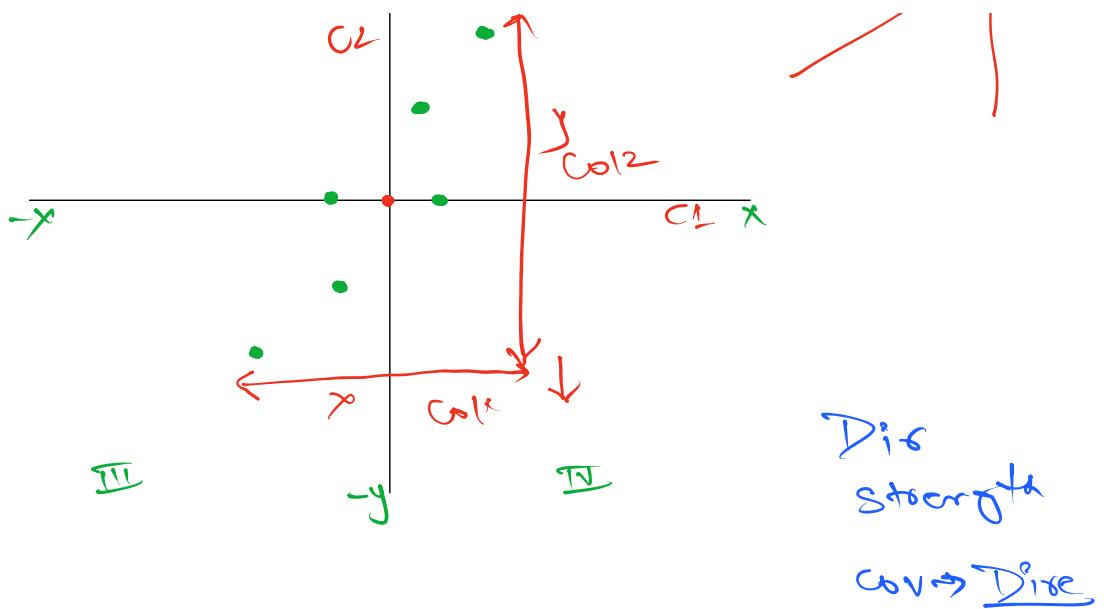
var(c)	Col1	Col2
$x_1 - 3$	-4	y_1
$x_2 - 1$	-2	y_2
$x_3 - 1$	0	y_3
$x_4 - 1$	0	y_4
$x_5 - 1$	2	y_5
$x_6 - 3$	4	y_6
	(0, 0)	

II

I

I





② Calculate the covariance matrix.

	Col1	Col2
Col1	$\text{var}(c_1)$	$\text{cov}(c_1, c_2)$
Col2	$\text{cov}(c_2, c_1)$	$\text{var}(c_2)$

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

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

$$\text{var}(c_1) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = (-3-0)^2 + (-1-0)^2 + (-1-0)^2 + (1-0)^2 + (1-0)^2 + \underbrace{(3-0)^2}_{6-1}$$

$$= 4.4$$

By similar calculation,

$$\text{Var}(c_2) = 8.0$$

Now,

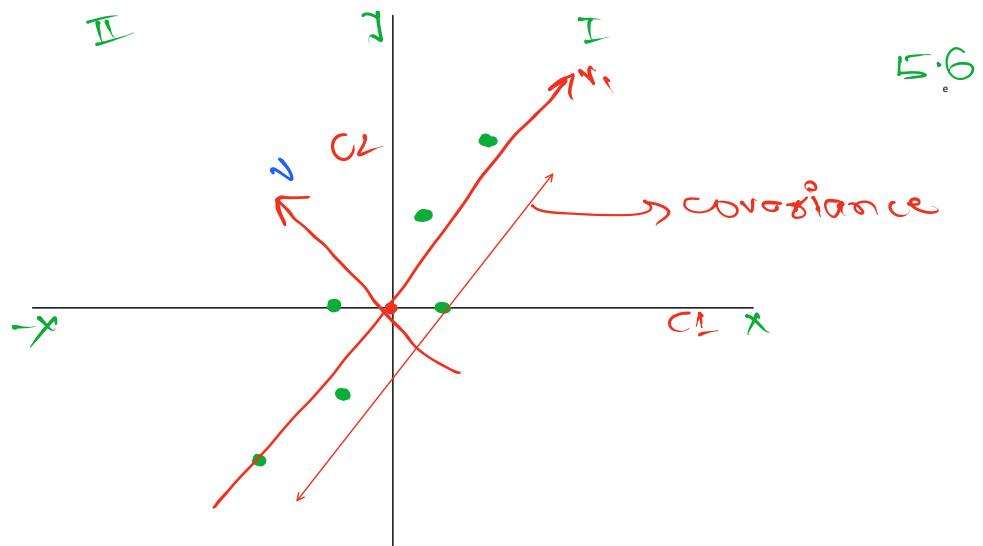
$$\begin{aligned} \text{Cov}(c_1, c_2) &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = (-3 \times -4) + (-1 \times -2) + \\ &\quad -1 \times 0 + 1 \times 0 + 1 \times 2 + 3 \times 4 \\ &= 6-1 \end{aligned}$$

$$\text{Cov}(c_1, c_2) = 5.6$$

$$\text{Cov Matrix} = \begin{bmatrix} 4.4 & 5.6 \\ 5.6 & 8.0 \end{bmatrix}$$

As we see from the covariance Matrix, that variance of Col2 is more than Col1.

And the covariance value is somewhere between these two variances



III -y | IV

- ③ Now we will calculate the eigenvalue and eigen vectors of covariance matrix. And for eigenvector we use this formula → Cov Matrix

2

$\begin{bmatrix} 2 & 4 \\ 2 & 8 \end{bmatrix}$

$$\rightarrow \det(A - \lambda I) = 0$$

$$\begin{bmatrix} 2 & 4 \\ 6 & 8 \end{bmatrix}$$

$$\det \left(\begin{bmatrix} 4.4 & 5.6 \\ 5.6 & 8.0 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \right) = 0$$

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Here the identity matrix should have same shape as covariance matrix.

$$\det \begin{bmatrix} (4.4 - \lambda) & 5.6 \\ 5.6 & (8.0 - \lambda) \end{bmatrix} = 0$$

$$(4.4 - \lambda)(8.0 - \lambda) - 5.6 \times 5.6 = 0$$

$$3 \cdot 84 - 12 \cdot 4 \lambda + \lambda^2 = 0$$

$$f(\lambda)$$