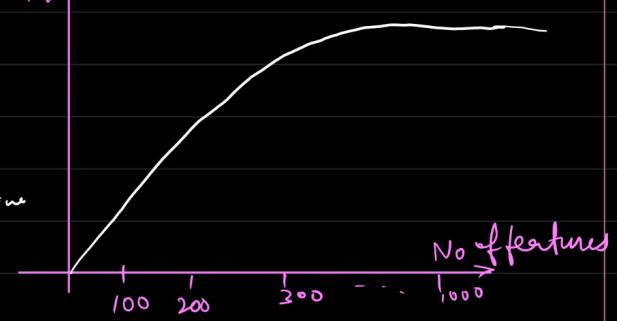


PCA (Principal Component analysis)

Requr ✓

20 f	50 f	100 f	250	500	1000
M ₁	M ₂	M ₃	M ₄	M ₅	M ₆

Acc / R² < Acc / R_{que} < Acc / R_c < Acc / R_{sg} \approx Acc / R_{sg} \approx Acc / R_{que}



* With increase in no of features, after one point of the Acc / R² stops increasing.

Why ??

→ Multicollinearity ($f_1, f_2, f_3 \approx f_4$)

→ All the entries of features are exactly same.

f_1

$$\left\{ \begin{array}{|c|} \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline \end{array} \right\}$$

→ two features are same

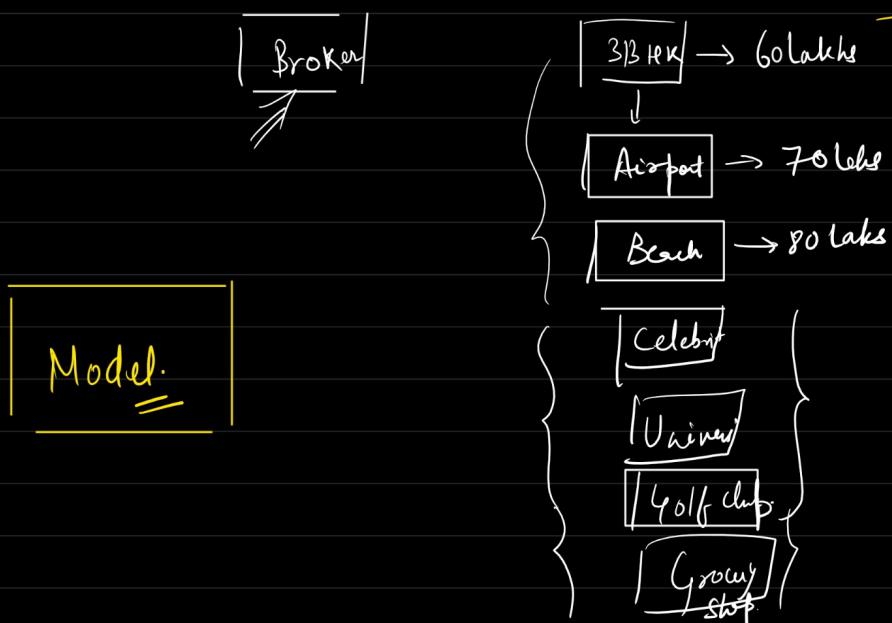
$$\begin{matrix} f_1 & f_2 \\ | & | \\ | & | \\ | & | \end{matrix}$$

→ lots of duplicate entries. f_1

→ No variance

1	1
1.01	1
1.02	2
1.03	2
1.04	2
= 1	
-	
-	

* With increase in no of feature performance of model degrades.



* Curse of dimensionality :-
With increase in no of feature
the performance of model degrades.

* To remove curse of Dimensionality :-

① Feature Selection

→ Correlation, VIF,

P-value, Select K feature,

I.G., Covariance, RFE

Disadvantage

→ dropping data

→ Time consuming.

② Feature Extraction



PCA (Dimensionality reduction technique)

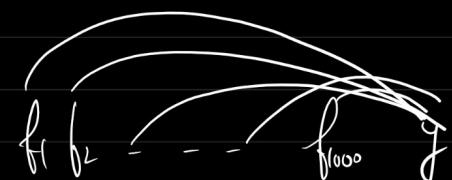
Why we should remove C.O.D?

- ① To improve performance of model.
- ② Visualise the data \rightarrow for insights
- ③ Prevents from overfitting
- ④ Better interpretation
- ⑤ To remove C.O.D

$$\frac{f_1 \ f_2 \ \dots \ f_{1000}}{\substack{\uparrow \\ 1000d \\ \downarrow \\ 3d, 21}}$$

* feature selection
 \rightarrow corr, cov.

	Area of room	No. of rooms	Price of house
1	100	1	1000
2	200	2	2000
3	300	3	3000
4	400	4	4000
5	500	5	5000
6	600	6	6000
7	700	7	7000
8	800	8	8000
9	900	9	9000
10	1000	10	10000



* Feature Extraction

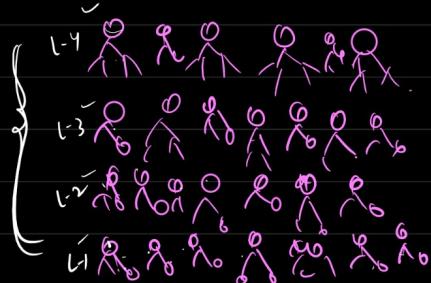
↓
 data transformation to extract
 new feature which represents
 both the feature.

(Domain understanding)	House Area	Price
-	-	-
-	-	-
-	-	-
-	-	-
-	-	-
-	-	-
-	-	-
-	-	-

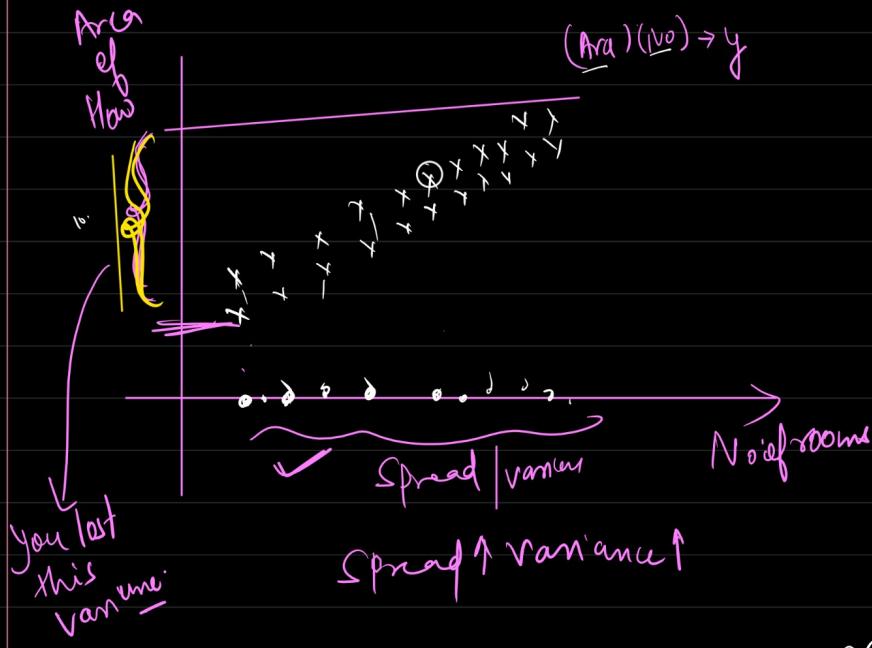
but say # of balloon	dictary provider	Distance from university
-	-	-

$$\frac{f_1 \ f_2 \ \dots \ f_{1000}}{\substack{\uparrow \\ \downarrow \\ \dots}}$$

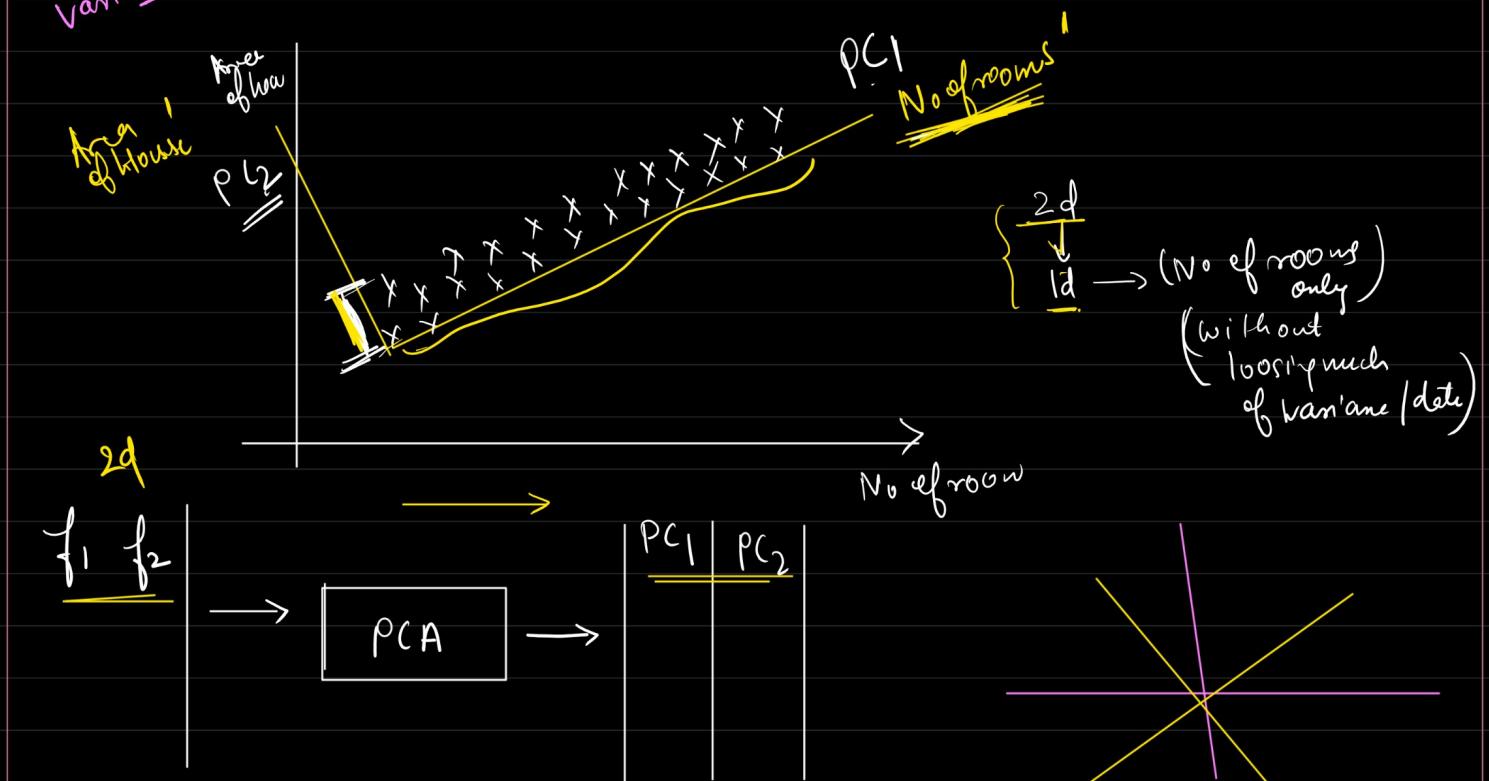
* PCA



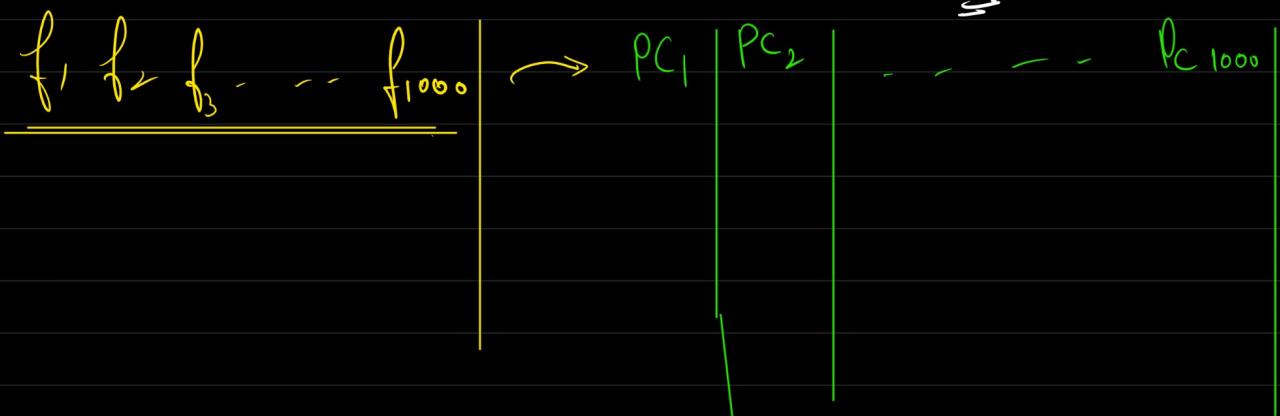
\rightarrow Variance / spread is more



$$(\text{Area})(\text{No}) \rightarrow y$$

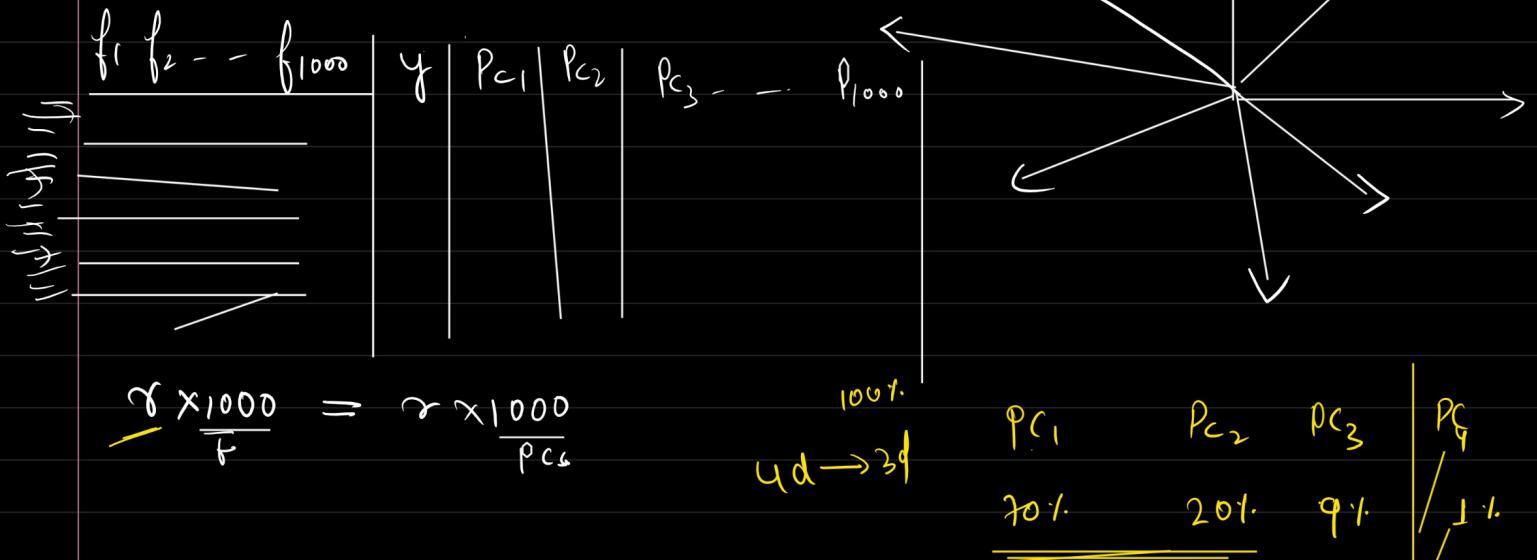


As PC_1 has maximum variance,
you will select only PC_1 for
model training.

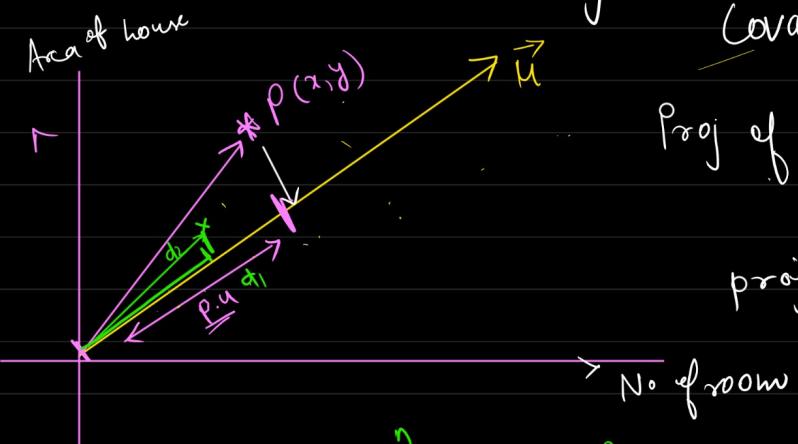


Two characteristics

- ① All of these PC's will be \perp to each other.
- ② \checkmark $PC_1 > PC_2 > PC_3 > PC_4 = \dots > PC_n$
- ③ No of PCs = No of features



* Axis transformation \Rightarrow Eigen decomposition of Covariance matrix.



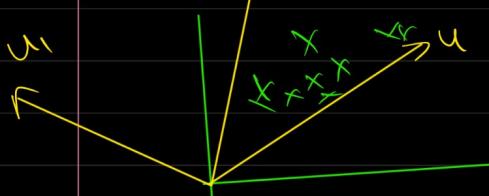
$$\text{Proj of } p \text{ on } u = \frac{p \cdot u}{\|u\|} \rightarrow \perp$$

$$\text{proj of } p \text{ on } u = p \cdot u$$

$$\left| d_1, d_2, d_3, \dots, d_n \right|$$

You want that Unit vector where Variance spread is maximum

Aim: To find that unit vector which captures maximum variance after projection



Covariance matrix

$$\overline{\text{Var}} = \frac{\sum (x_i - \bar{x})^2}{n}$$

$$\text{Covariance} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

	x_1	x_2
x_1	$\text{Cov}(x_1, x_1)$	$\text{Cov}(x_1, x_2)$
x_2	$\text{Cov}(x_2, x_1)$	$\text{Cov}(x_2, x_2)$

$$\begin{bmatrix} \text{Var } x_1 & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_1, x_2) & \text{Var } x_2 \end{bmatrix}$$

$$\begin{array}{cccc} & x_1 & x_2 & x_3 \\ \begin{array}{|l|l|l|} \hline x_1 & \text{Var}(x_1) & \text{Cov}(x_1, x_2) & \text{Cov}(x_1, x_3) \\ \hline x_2 & \text{Cov}(x_2, x_1) & \text{Var}(x_2) & \text{Cov}(x_2, x_3) \\ \hline x_3 & \text{Cov}(x_3, x_1) & \text{Cov}(x_3, x_2) & \text{Var}(x_3) \\ \hline \end{array} & \end{array}$$

Covariance matrix
↔
Variance / Covariance matrix.

whole idea
is to
capture spread.

Theorem: → If you decompose a covariance matrix of features, then you will get eigen value and Eigen vector and eigen vector with highest magnitude of eigen value capture the maximum variance / spread

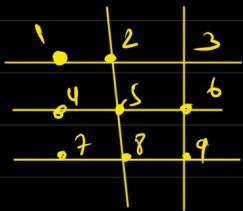
- 2nd highest magnitude → 2nd highest var
- 3rd highest variance → 3rd highest var
- 4th → 4th highest.

$$P(C_1 > P(C_2) > \underline{P(C_3)}$$

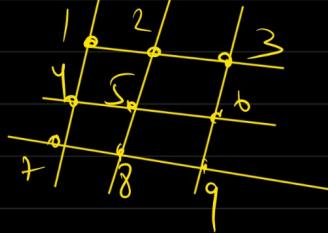
* Linear transformation of a matrix.

Analogy

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

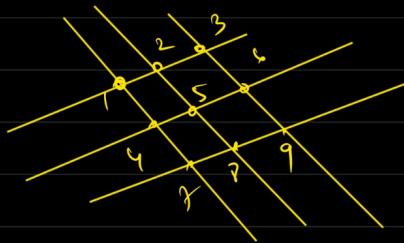


linear transform



Linear transformation

↳ A matrix transformation
that brings changes
in the coordinate
system.



$$A \vec{v} = \lambda \cdot \vec{v}$$

↓ ↓ ↗
 matrix Eigenvalue Eigenvalue Eigen vector

When we use covariance matrix
as A, Eigen decomposition
of covariance matrix

Eigen Value Eigen Vector

Linear transformation pointers

→ Many vectors were changing both
magnitude & direction.

→ Few vectors changed only
magnitude & not the
direction.

$$\begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix}$$

$$\parallel (1, 0) \rightarrow (3, 0) \text{ direction } \begin{matrix} \text{same, magnitude} \\ \text{changed} \end{matrix}$$

$$\parallel (-1, 1) \rightarrow (-4, -2) \text{ direction } \begin{matrix} \text{same, magnitude} \\ \text{changed} \end{matrix}$$

These vectors are called
as Eigen vector.

* Eigen values are
change in magnitude
for eigen vector

$$A \vec{v} = \lambda \vec{v}$$

↓ ↗
 Covariance matrix Scaler

* The vector which only changes magnitude and not direction will be equals to dimension of matrix / no of feature.

$$\text{No of feature} = \text{No of PC}$$

$$A \underline{v} = \lambda \underline{v} \rightarrow \text{Vector} \\ \downarrow \text{stretch} \quad | \quad \text{shrink}$$

Steps to calculate Eigen value / Eigen vector (PC's)

Step-1 Standardise the date (make the date mean centred)

Step-2 Covariance matrix ($\text{df} \cdot \text{cov}()$)

Step-3 Eigen decomposition np. linalg

$$A \underline{v} = \lambda \underline{v} \\ \lambda \rightarrow \text{eigen value} \\ \underline{v} \rightarrow \text{eigen Vector}$$

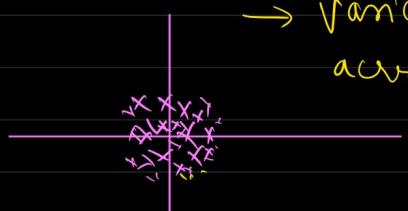
$$1000 \left\{ \rightarrow 1000 \text{ PC's}$$

$$\text{PC}_1 > \text{PC}_2 > \dots > \text{PC}_n$$

* Max var will be captured by first few principal comp.

* PC's are \perp^r to each other

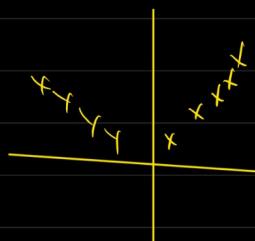
* Disadvantage of PCA



→ Variance spread across all axis is

Same.

→ PCA will fail

 → PCA will
lose the
pattern.



* MNIST

0 — φ

 → Pixels → 28×28