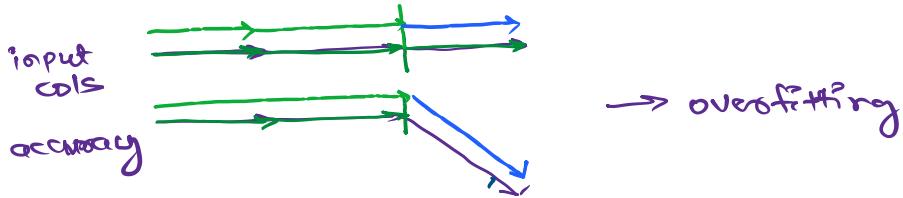


Age	Floors	Rooms	Area	Price

Curse of Dimensionality



When you have too many input cols, then there are high chances that some/many of these cols are unnecessary, because of this we start to see overfitting in the data. This problem is what we refer as Curse of Dimensionality.

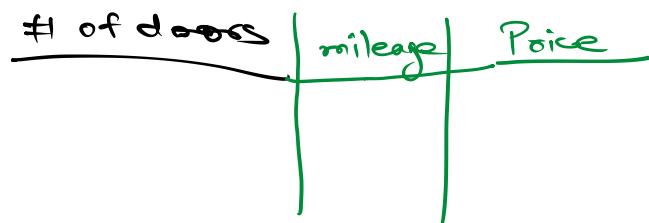


There are also high chances of multi-collinearity.

↓
High corr. b/w
input cols

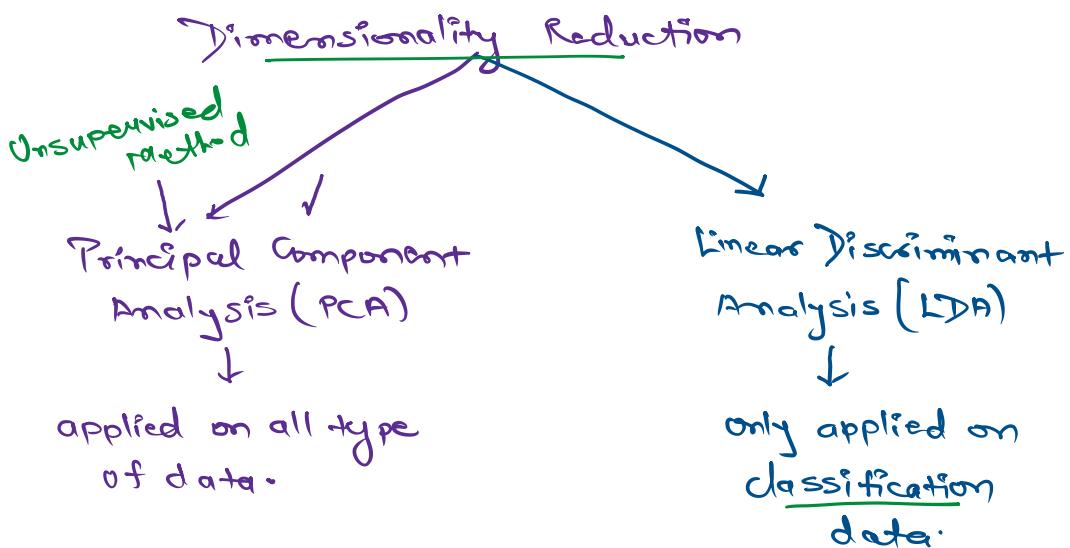
# rooms	# of floors	Area	L	W	Greyson	Price

Dataset \rightarrow 200 cols 200D
 \downarrow
 28 cols \rightarrow unnecessary \rightarrow Delete them
 \downarrow
 172 remaining cols
 \downarrow
 56 cols are very less important



Solution: If we have some method which can reduce the number of input cols in a smarter way, rather than directly deleting it, that would be a very good solution.

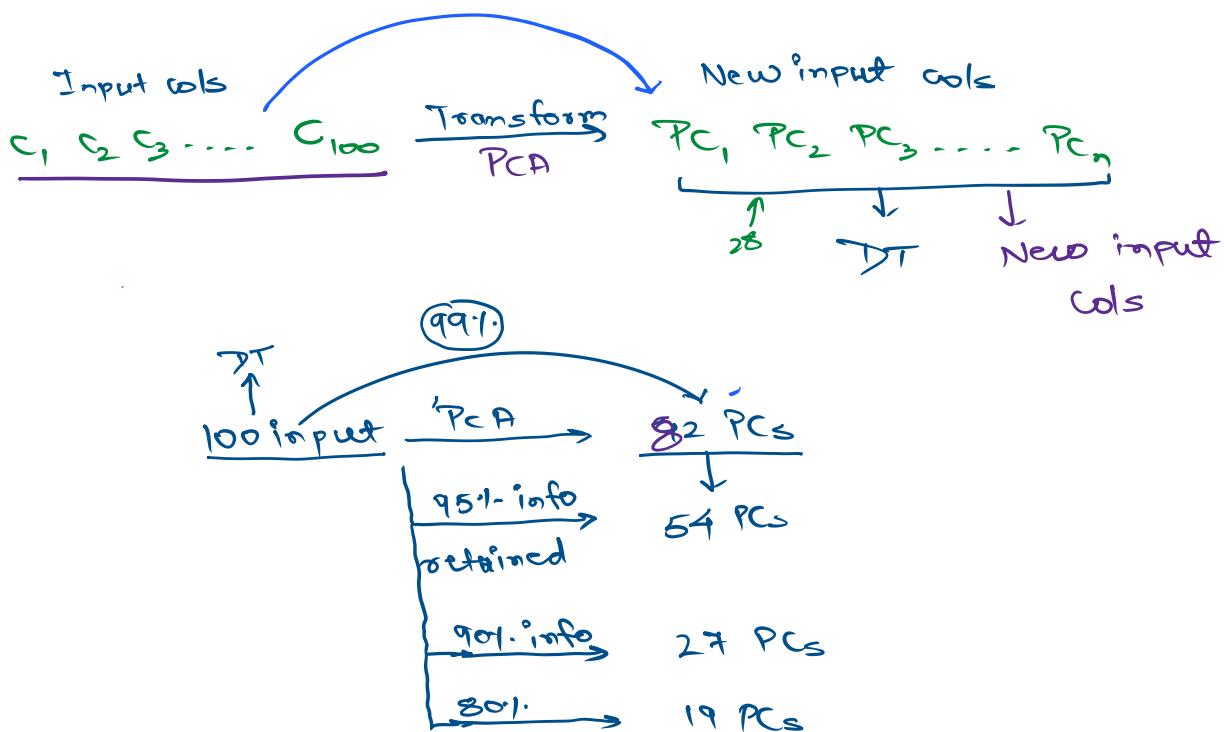
\downarrow
 This is where 'Dimensionality Reduction' technique come into picture.



Principal Component Analysis

PCA helps reduce the dimension of the data in such a way that we lose very little information while reducing the no. of input cols.

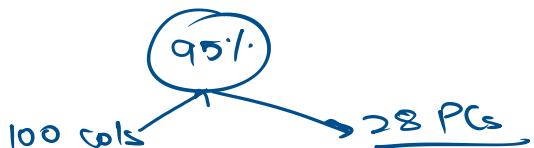
$$2.5C_1 + \frac{1}{17}C_{22} + 3C_3$$



information contained in PCs decrease from PC_1 to PC_{28}

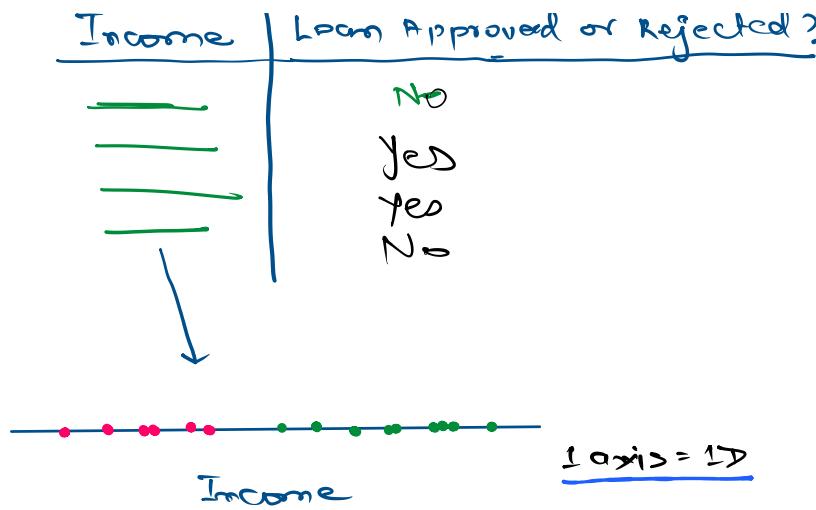
$PC_1, PC_2, PC_3, PC_4, \dots, PC_{28}$

Combination of
few most imp.
input cols.
(27%)

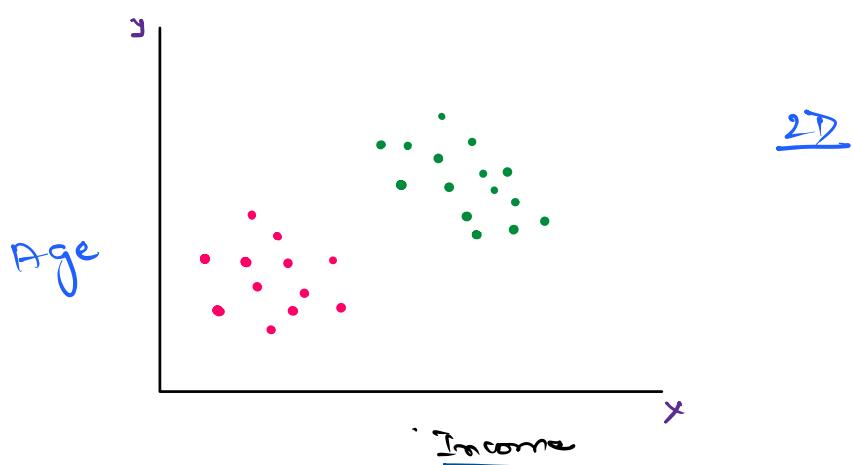


Mathematical Intuition of PCA \rightarrow

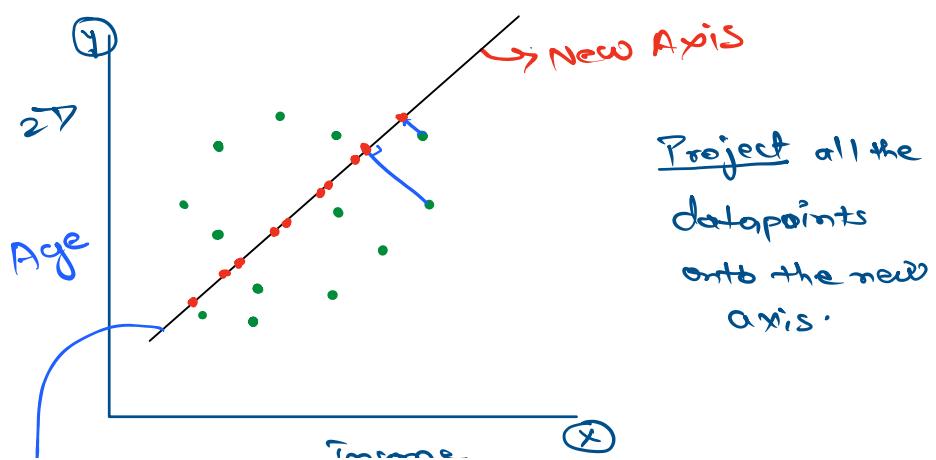
$$\text{Inp} \downarrow \quad | \quad \text{Out} \downarrow$$

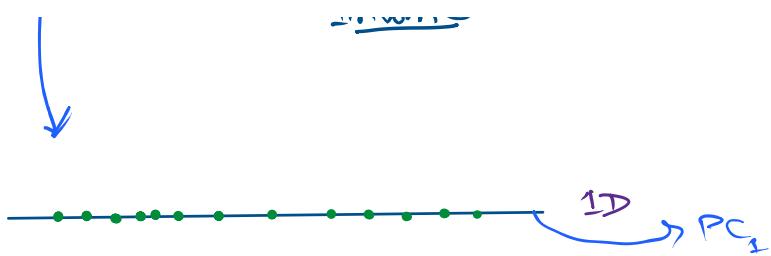


Income	Age	Loan App or Rej?
—	—	—



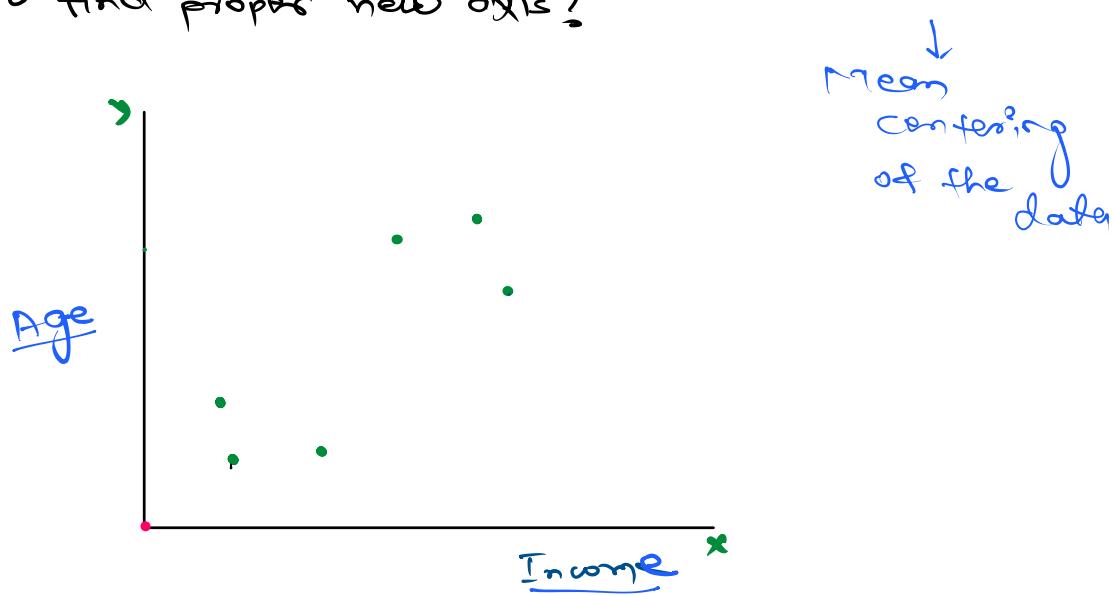
How to reduce the dimension using PCA →



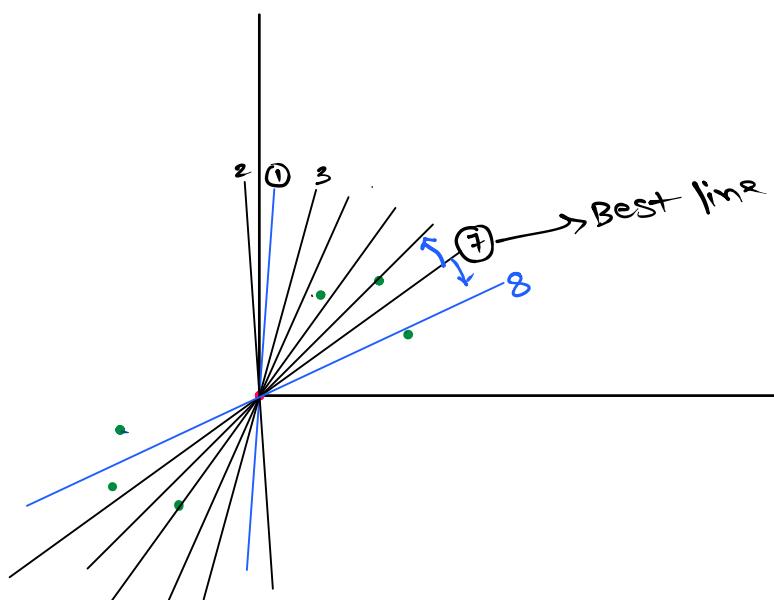


Projection of datapoints onto the new axis is nothing but linear combination of existing cols/dimension.

How to find proper new axis?

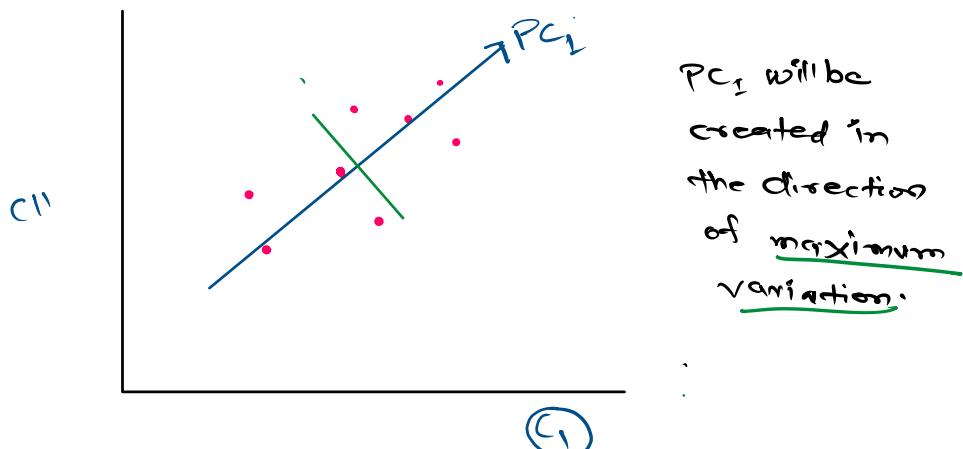
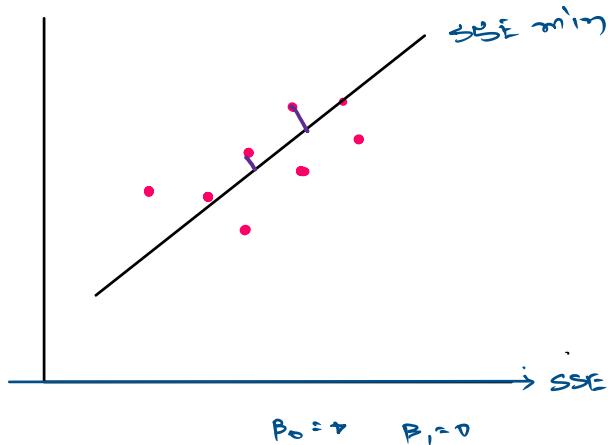


We have to move the data in such a way that its center is at 'origin'.



$$d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2 \quad (\text{SSE})$$

Any line which minimizes this sum will be our new axis.



mileage	# of doors	Price
17	4	
23	2	
18.6	4	
12	4	
13	2	
30	4	

Standardization \rightarrow

$$\frac{x - \bar{x}}{\text{std}}$$

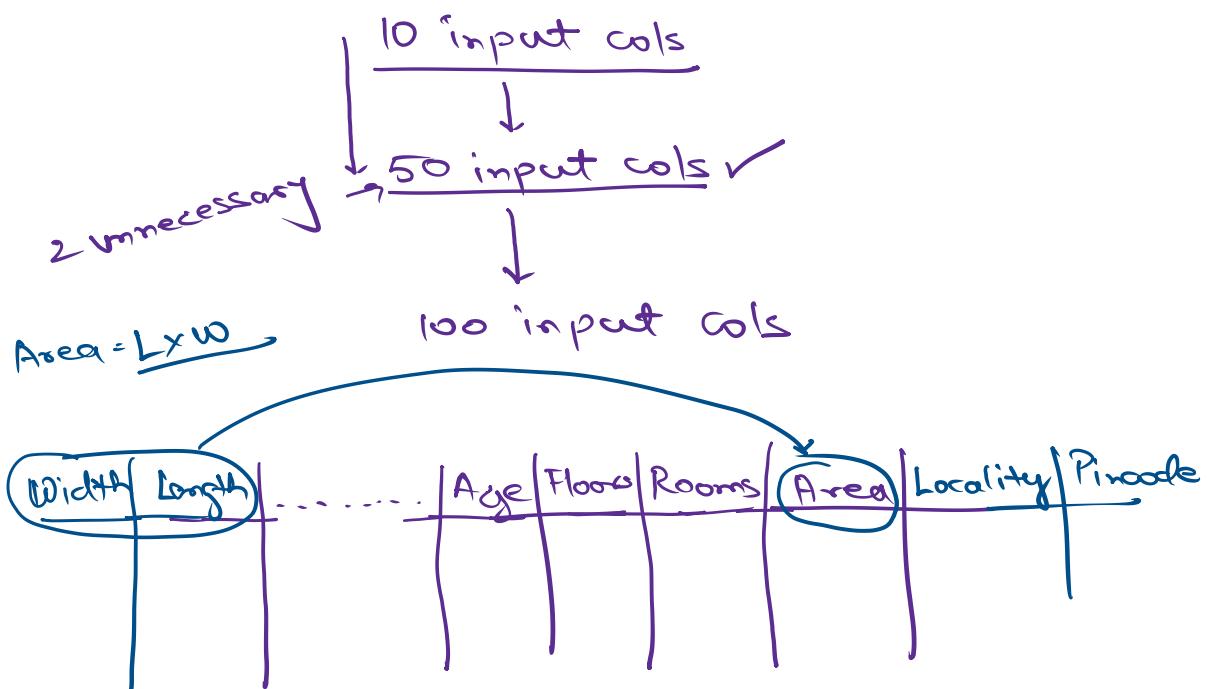
Linear Regression

$$x_1 = \# \text{ of rooms} \rightarrow 4$$

$$x_2 = \# \text{ of floors} \rightarrow 2$$

$$x_3 = \text{Area of house.} \rightarrow 2000$$

$$\begin{aligned} J_{\text{pred}} &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 \\ &= 1.2 + 0.9 x_1 + 0.3 x_2 + 2.1 x_3 \\ &= 1.2 + 0.9 \times 4 + 0.3 \times 2 + 2.1 \times 0.2 \\ &= > 4200 \text{ error.} \end{aligned}$$

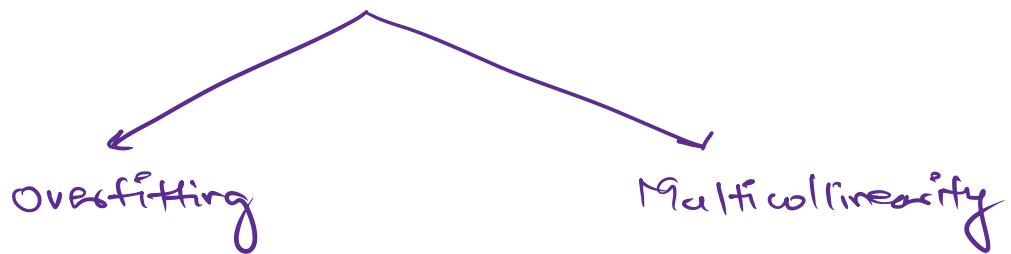


Curse of Dimensionality

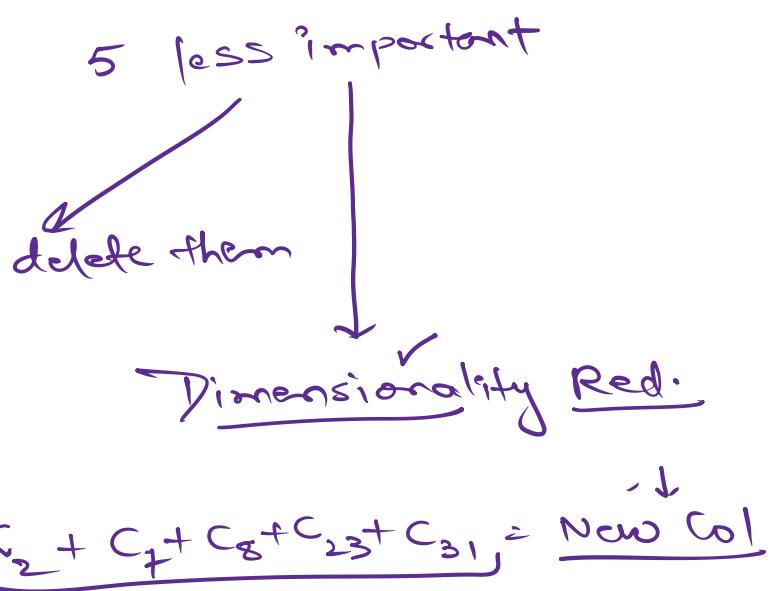
Unnecessary & redundant cols ✓



These will create noisy
data



Less important



→ PCA → unsupervised

