

Principle of Component Analysis (PCA)

* Curse of dimensionality
 \Downarrow
 features.

To predict price of house
 $f_1 \ f_2 \ \dots \ f_{100} \mid y$



$Acc_{Rsquare} < Acc_{Rsquare} \uparrow < Acc_{Rsq} \uparrow < Acc_{Rsq} \uparrow \approx Acc_{Rsq} \approx Acc_{Rsq} \approx Acc_{Rsq}$

* With increase in no of features,
 after one point of time the
 $Acc(Rsquare)$ will not increase
 in the proportion.

\Downarrow
 why?

→ Few of features will be multicollinear

($f_1, f_2, f_3 \approx f_4$)
 → few of features might be exactly
 same

$f_1 = 1, 1, 1, 1, 1, 1, 1, 1$
 $f_2 = 1, 1, 1, 1, 1, 1, 1, 1$

→ No variance / information in feature.

→ lots of duplicate entries

f_1
 1.01
 1.02
 1

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

Analogy

* you want to buy a house.

$\boxed{\text{Brokers}}$

$\boxed{2 \text{ BHK}} \rightarrow 60 \text{ Lakhs}$

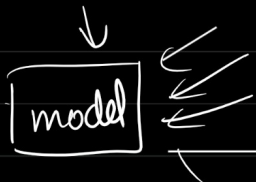
$\boxed{3 \text{ BHK}} \rightarrow 80 \text{ Lakhs}$

$\boxed{\text{beach}} \rightarrow \uparrow \uparrow$

$\boxed{\text{Airport}} \rightarrow \uparrow \uparrow$

Execution time

Evaluation metrics



near celebrity
Grocery
Shopping
University

* Curse of dimensionality \rightarrow With increase in no of feature the performance of model degrades.

To remove COD :-

- ① Feature Selection ② Feature Extraction

\downarrow
 PCA (Dimensionality reduction technique)

* Why we should remove COD?

- ① To improve the performance of Model.
- ② Visualise the data \rightarrow for insights
- ③ Prevents from overfitting.
- ④ Better interpretation.
- ⑤ To remove Curse of dimensionality

1000d
 \downarrow
3d

* Feature Selection

Area of house (x_1)	Near to airport (x_2)	Price of house. y

① Correlation (Pearson)

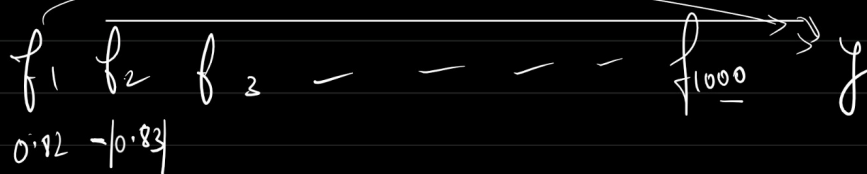
$$\text{Corr} = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

$x \uparrow y \uparrow$
 $x \downarrow y \downarrow$
 $x \uparrow y \downarrow$
 $x \downarrow y \uparrow$

$$\text{② Cov}(x, y) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

$$\text{Corr}(\overset{\downarrow}{x_1}, \overset{\downarrow}{y}) = 0.82 \checkmark$$

$$\text{Cor}(x_2, y) = \underline{\underline{0.3}}$$



* Feature Extraction

(In feature selection, you dropped the features. But let's if you want all the features and also curse of dimensionality, then comes the technique feature extraction)

[illegible]

data transformation to extract New feature which represents both the features.

House Area	Price
using (Domain expertise)	-

but say \rightarrow # of balloons | distance from airport | distance from University.

Principal component Analysis

→ Dimensionality reduction technique.