

Curse of Dimensionality

1) binary features $\Rightarrow f_1^{(3-d)} f_2 f_3 = \# \text{ datapoints} = 2^3$

10-d $\Rightarrow \# \text{ datapoints} \Rightarrow 2^{10}$

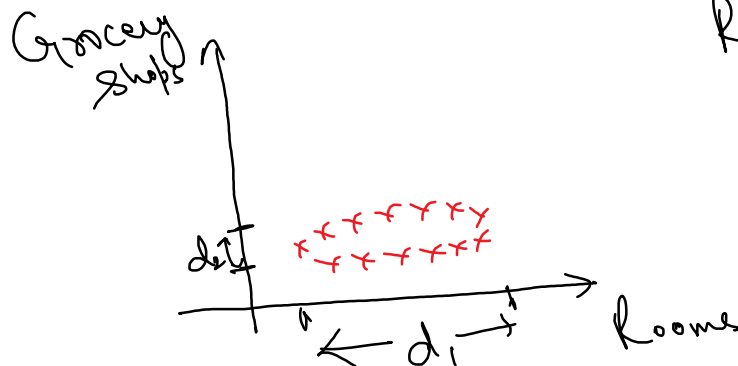
$\# \text{ dimensions} \uparrow \propto \text{overfitting} \uparrow \rightarrow \text{model performance}$

2) Distance f'ns \rightarrow Euclidean distance

NLP \Rightarrow Hamming / cosine similarity

Feature Extraction (PCA)

Feature Selection



Rooms Grocery shops Price of flat

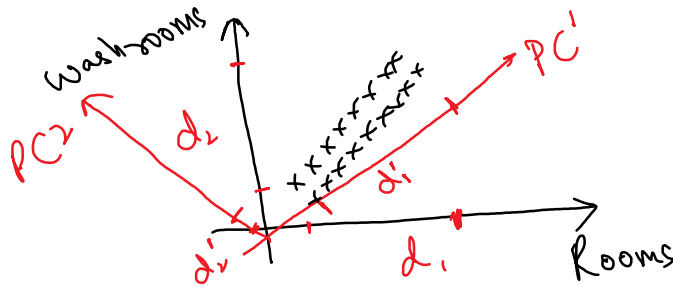
$d_1 > d_2$

Rooms will be chosen

higher variance \rightarrow higher info

Grocery Shops will be dropped

Feature Extraction



$d_1 \approx d_2$ $d_1' > d_2'$

Rooms & Washrooms
 \downarrow
Size

Projection of \vec{x} on \vec{u} (unit) $= \frac{\vec{u} \cdot \vec{x}}{\|\vec{u}\|} = u^T x \rightarrow [LA]$

MSE $= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$ $\Rightarrow \sum_{i=1}^n \frac{(u^T x_i - u^T \bar{x})^2}{n}$

mean \downarrow up hard

PCA

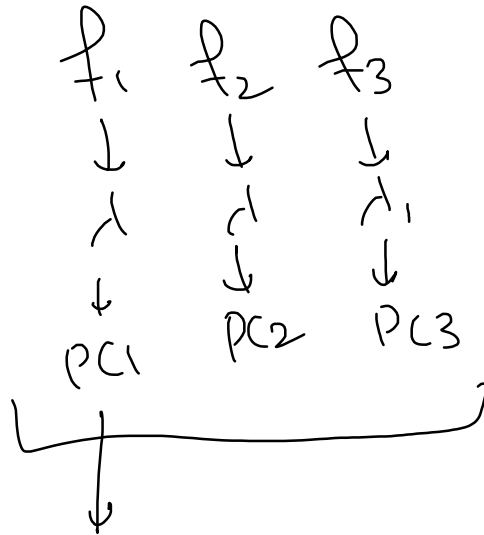
1 \rightarrow Mean Centering

2 \rightarrow Co-variance Matrix $\begin{matrix} \lambda_1 & \lambda_2 & \lambda_3 \end{matrix}$

$$\begin{matrix} f_1 \\ f_2 \\ f_3 \end{matrix} \begin{bmatrix} \text{var } f_1 & \text{cov}(f_1 f_2) & \text{cov}(f_1 f_3) \\ \text{cov}(f_1 f_2) & \text{var } f_2 & \text{cov}(f_2 f_3) \\ \text{cov}(f_1 f_3) & \text{cov}(f_2 f_3) & \text{var } f_3 \end{bmatrix}$$

3 → Eigen decomposition

Info comparison
 $PC1 > PC2 > PC3$



$n\text{-components} = 1$

→ 11

⇒ PC1