# Curse of Dimensionality

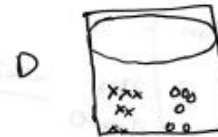$\downarrow$ $d$ is high
'dimension'

## ① ML   Binary features

$\begin{cases} f_1, f_2, f_3 \longrightarrow \text{Total \# datapoints} = 2^3 = 8 \\ \\ \qquad\qquad\qquad\qquad\qquad = 2^{10} = 1024 \\ \\ f_1, f_2, f_3, \ldots, f_{10} \longrightarrow \end{cases}$

as dim $\uparrow$; the # datapoints to
perform good model
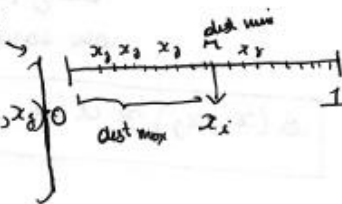increase exponentially

$D$

### Hughes phenomenon

'$n$' $\rightarrow$ size of dataset is fixed
performance $\downarrow$ as dim $\uparrow$.

## ② Distance functions (Euclidean distance)

Curse of Dim :- intuition of distance in 3D
is not valid in high dimensional spaces.

1-D world $\rightarrow$ n random pt

distance minimum $= \min\limits_{x_j \neq x_i} \begin{cases} dist(x_i, x_j) \neq 0 \end{cases}$
$(x_i)$

dist min
$x_1\, x_2\, x_3 \quad x_r$

dist max   $x_i$   1

$\downarrow$
distance to nearest point from $x_i$ where $x_i \neq x_j$

$$\text{distance maximum}(x_j) = \max_{x_j \neq x_i}\left\{ \text{euc dist}(x_i, x_j) \right\}$$

(3D) $\quad \dfrac{\text{dist max}(x_i) - \text{dist min}(x_i)}{\text{dist min}(x_i)} > 0$

when $d = 1, 2, 3$

as dim $\uparrow$

$$\left[ \lim_{d \to \infty} \frac{\text{distmax}(x_i) - \text{distmin}(x_i)}{\text{distmin}(x_i)} \to 0 \right]$$

As dimensionality increases $\qquad \uparrow$ ratio approaches 0.

$\Downarrow$

$$\left( \text{distmax}(x_i) \approx \text{distmin}(x_i) \right)$$

high dimensional space if you take

— n random pts

$\boxed{x_i}$ $\quad$ distmax$(x_i) \approx$ distmin$(x_i)$

$\Downarrow$

every pair of pts
are equally dist. from each other

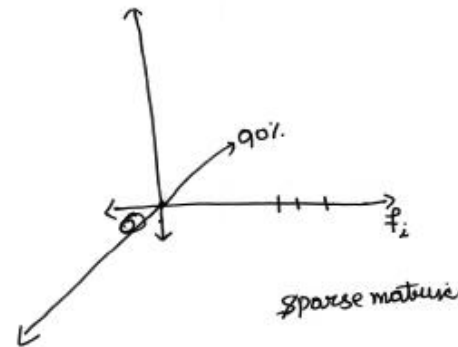$$\boxed{d(x_i, x_j) \approx d(x_i, x_k)}$$

3.) KNN → euc - distance
↓
high dimensional space

euclidean distance → logically not make sense.

( does not make sense in high dimension space.)

↳ solution : cosine similarity → high dimensionality
↳ but less effected

NOTE :-

Twist :- 1) If data is high dimensional & dense → impact of dim is high

ii) versus high dimensional & sparse → impact of dimensionality
↓  is lower.
not uniform random
/spread of data in the
sparse dimension.
(non zero less)



90%.

$f_i$

sparse matrix

③ Overfitting & underfitting

dimensionality ↑ , overfitting ↑ ⟶ linear regression.

Classification
oriented → forward feature selection : pick most useful subset
of features

→ Dimension - reduction :- PCA, tsne → do not use class
label and not classification oriented

→ KNN on text data

      ⌐→ cosine similarity instead of Euclidean dist.

      └→ sparse representation instead of dense representation

             ↓

         bag of words.