

Agenda

- Bias Variance of KNN
- Time & Space Complexity
- KNN Imputation

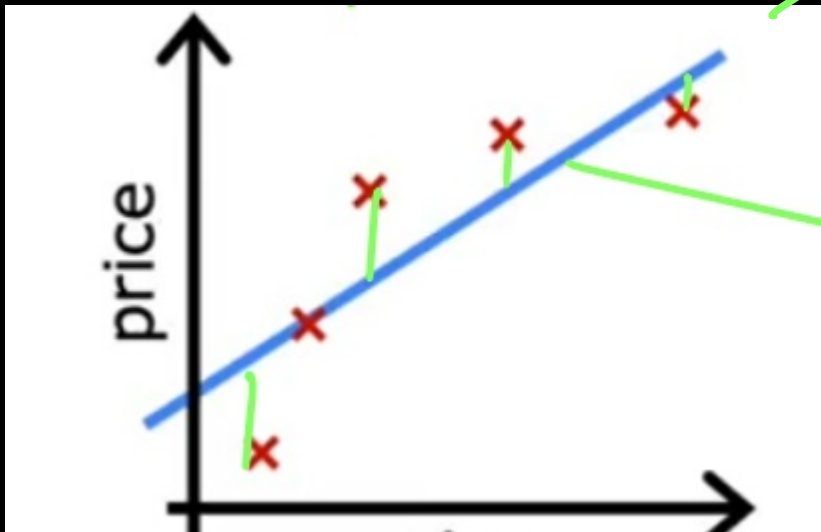
→ A : Underfitting
B : Overfitting

$$y = w_1 x + w_0$$

$$y = w_1 x + w_{12} x^2$$

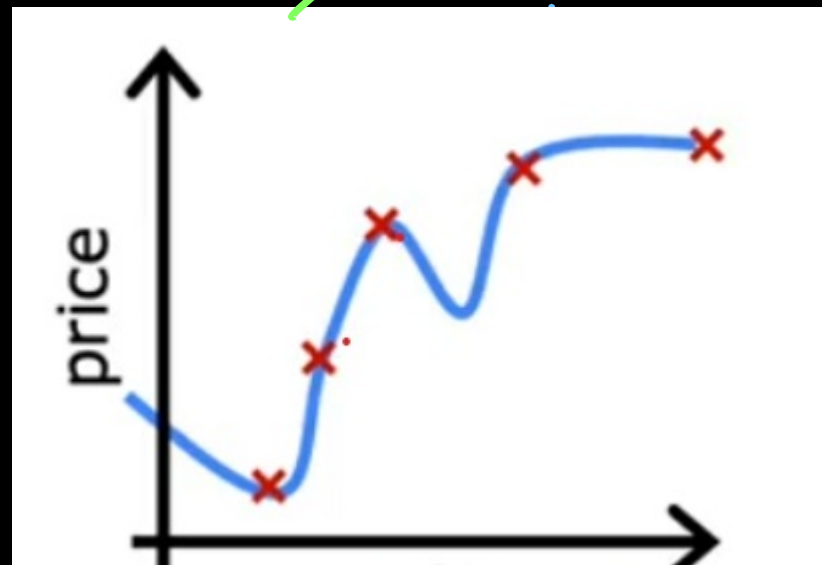
$$+ w_{13} x^3$$

$$+ w_{14} x^4 + w_0$$



(1)

→ High bias & low variance

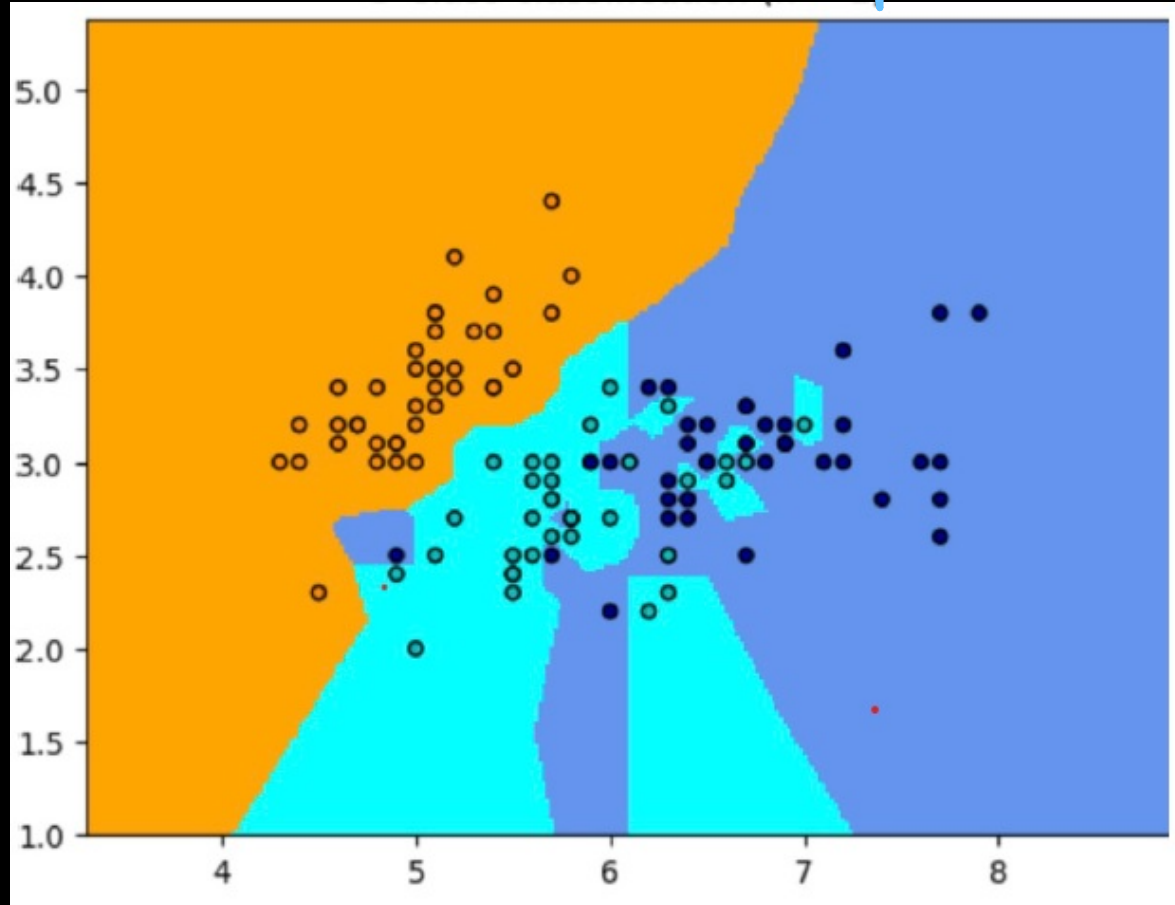


(2)

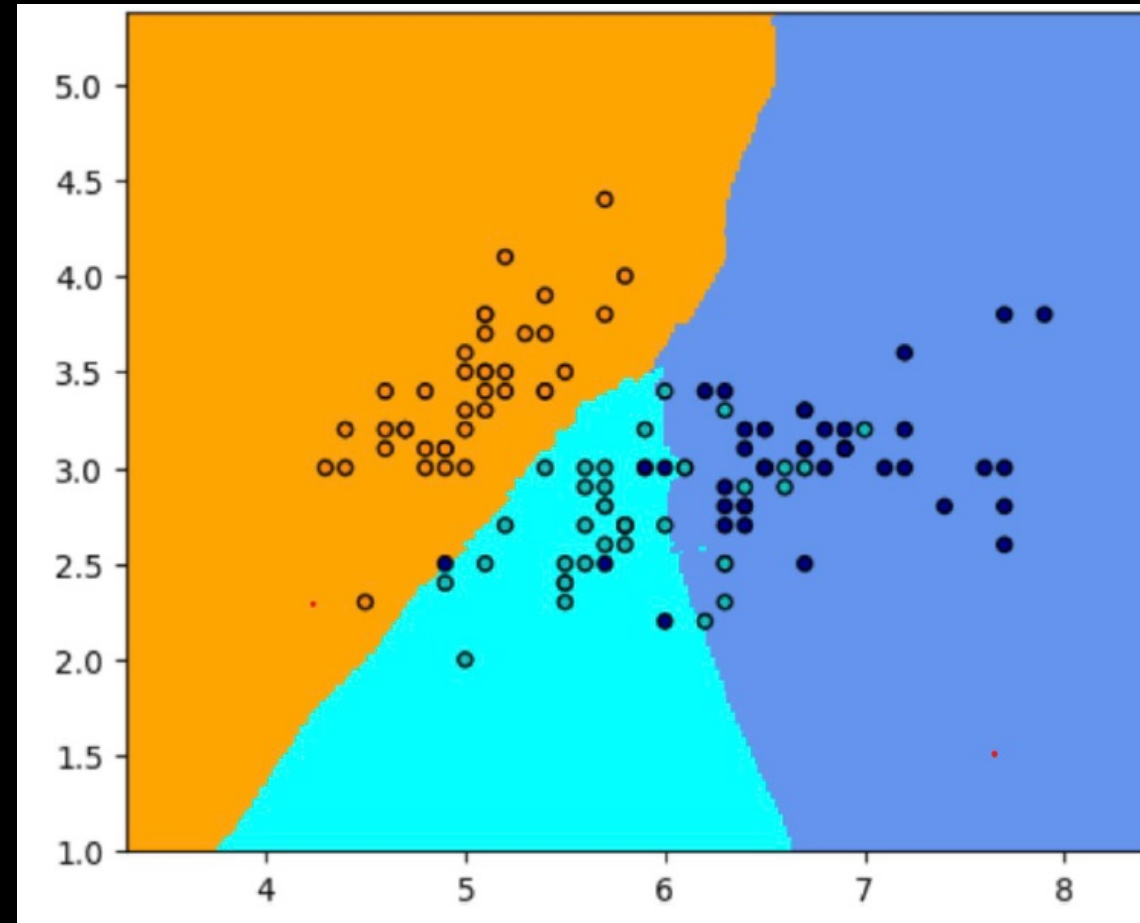
low bias

& High variance

✓ \rightarrow Decision boundary is more complicated



✓ \rightarrow Decision boundary is simple



(1) Overfitting

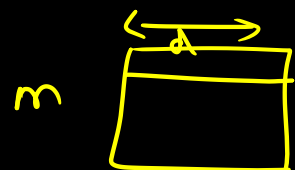
(2) Underfitting

$K = 1$

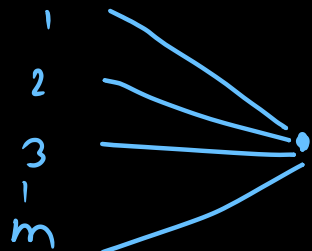
$K = 15$

K -NN

Algorithm



x_q



1. Compute the dist. of test point x_q from every training point $\rightarrow O(md)$
2. Take the top-K ^{smallest} closest distance $\rightarrow O(m \log m)$
3. For classification \rightarrow Find the majority with K class

$\hookrightarrow O(K)$

$\hookrightarrow O(md + m \log m + \cancel{K})$

$\hookrightarrow \boxed{O(md + m \log m)} \leftarrow$ For each test sample

$O(nmd + nm \log m)$

At test time

$$\hat{y}^{(i)} = \sigma(z^{(i)})$$

Single test
sample

$$= \sigma(w x^{(i)} + w_0)$$

$O(d)$

$$w_1 x_1^{(i)} + w_2 x_2^{(i)} - - - w_d x_d^{(i)} + w_0$$

← logistic

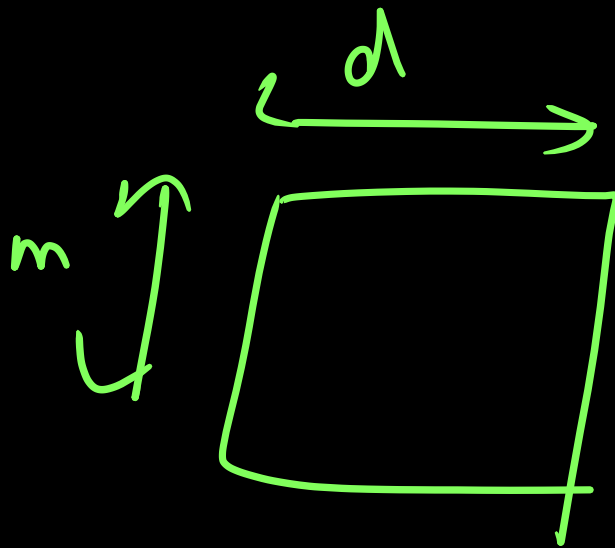
$$O(md + m \log m) \leftarrow \text{kNN}$$

1) Logistic — $O(dn)$

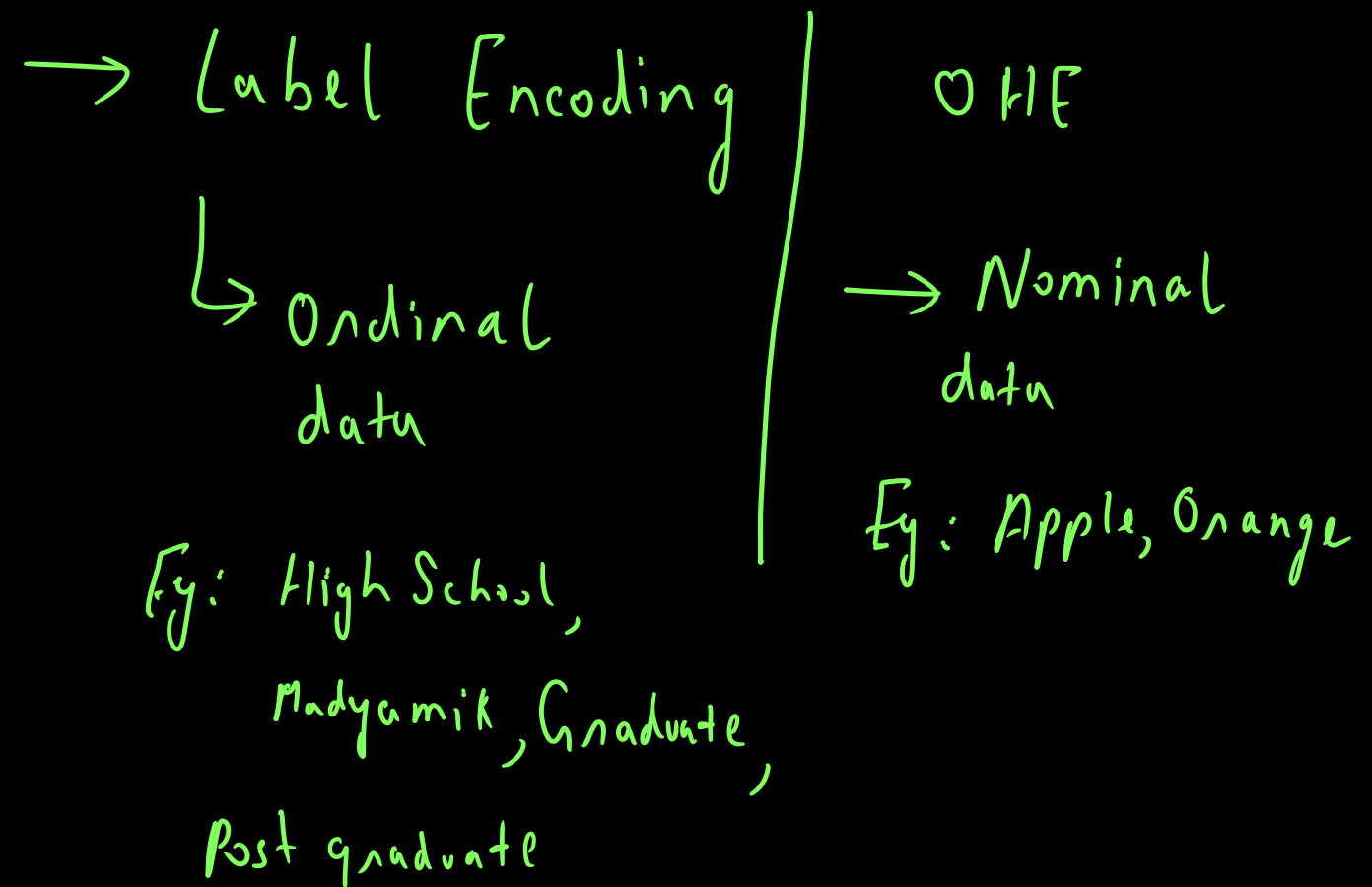
2) KNN — $O(\underbrace{md + m \log n}_{\checkmark}) \cdot n$

Space Complexity

↳ $O(md)$



Categorical data



↳ In 2 categories

M	0
F	1

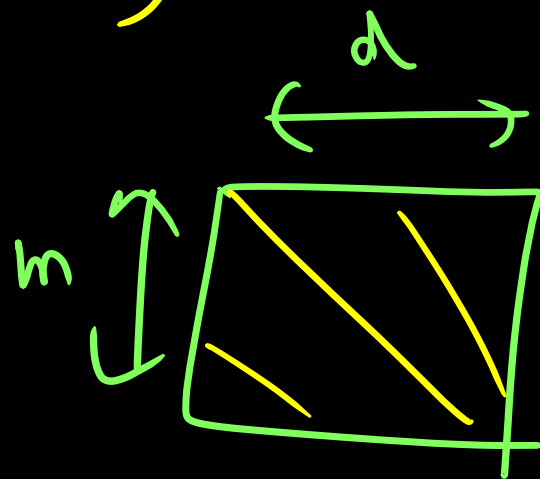
<u>Blood</u> <u>group</u>	<u>One hot</u>			
	A	B	AB	O
A	1	0	0	0
B	0	1	0	0
AB	0	0	1	0
O	0	0	0	1

Time & Space Complexity

Training time

Time $\rightarrow O(1)$

Space $\rightarrow O(md)$



Testing time

Distance

$x^{(1)}$ $x^{(2)}$

1) Manhattan (L_1) norm $\rightarrow \sum_{i=1}^d |x_i^{(1)} - x_i^{(2)}|$

2) Euclidean (L_2) norm

3) Minkowski

$$\rightarrow \left[\sum_{i=1}^d (x_i^{(1)} - x_i^{(2)})^2 \right]^{1/2}$$

4) Cosine Similarity

$$\rightarrow \left[\sum_{i=1}^d |x_i^{(1)} - x_i^{(2)}|^p \right]^{1/p}$$

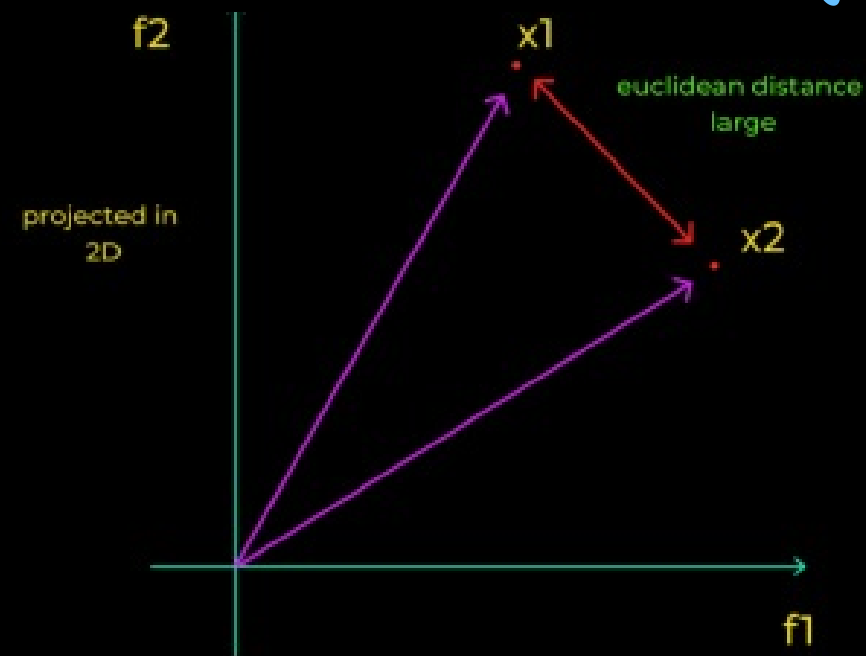
$p=1 \rightarrow \text{Manhattan}$

$p=2 \rightarrow \text{Euclidean}$

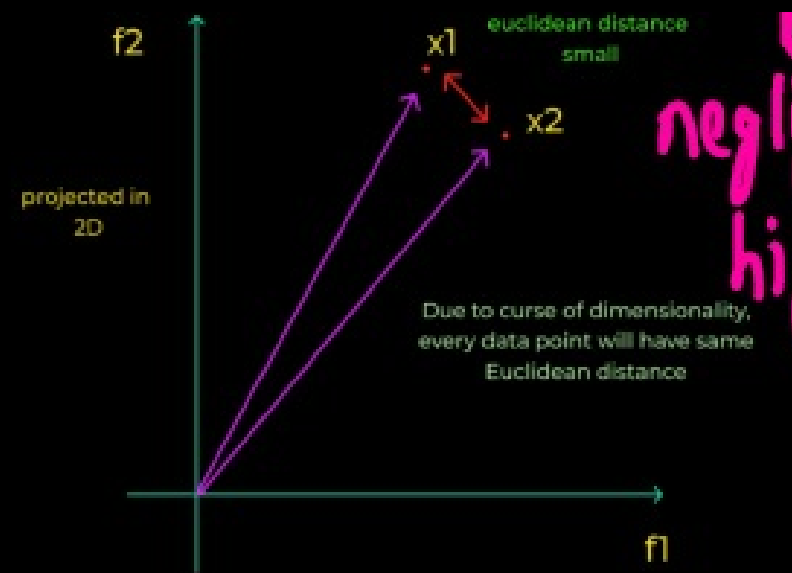
$$\text{Cosine Similarity } (x^{(1)}, x^{(2)}) = \frac{x^{(1)} \cdot x^{(2)}}{\|x^{(1)}\| \|x^{(2)}\|}$$

$\hookrightarrow (-1, 1)$

Euclidean \rightarrow Cannot be used with
high dimension data



Due to low dimension,
Euclidean distance
between $x1$ & $x2$ is very
large



negligible
high

Euclidean distance
cannot be used

Due to high dimension,
Euclidean distance
between $x1$ & $x2$ is very
small



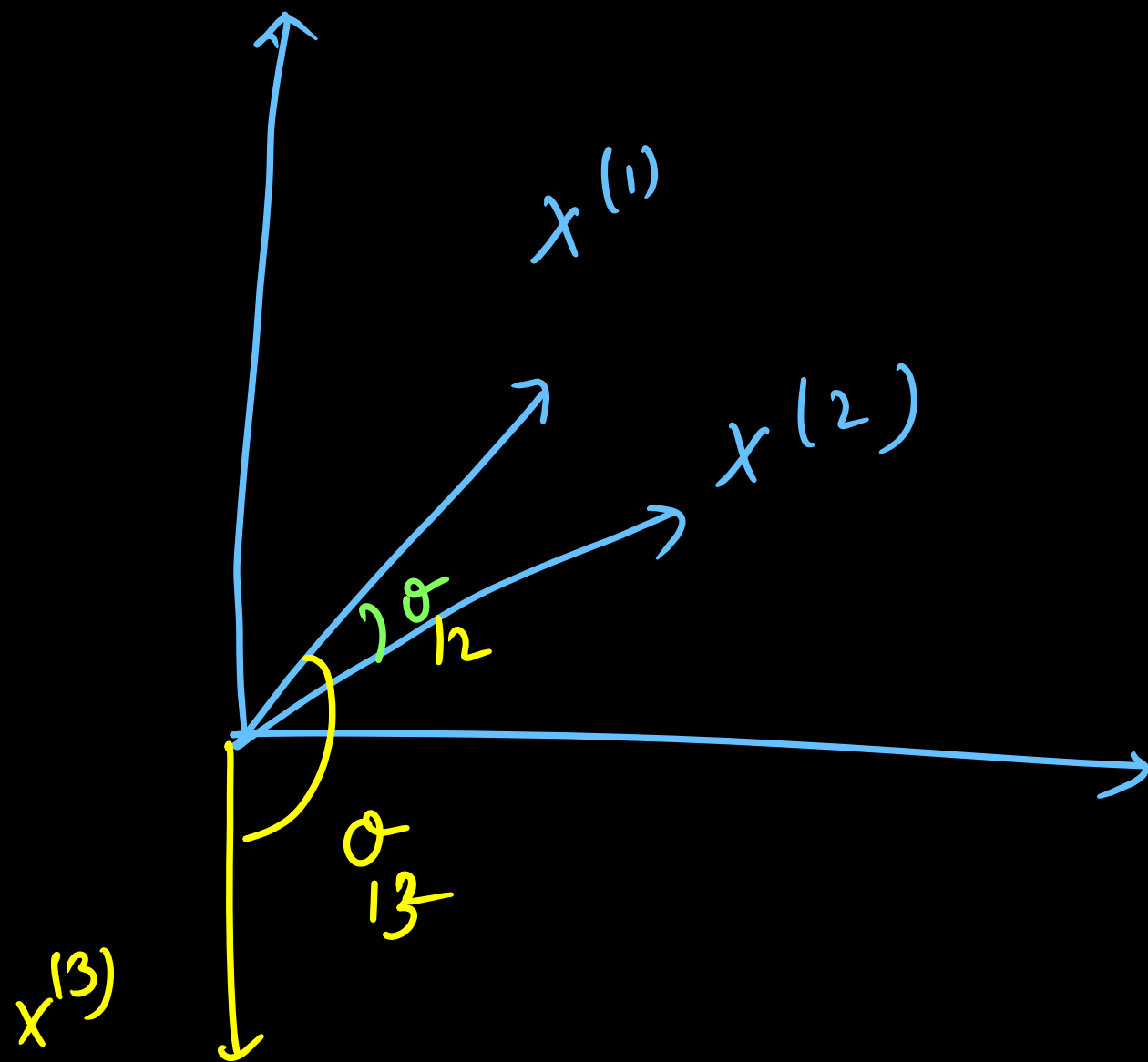
Conclusion : Euclidean distance fails when dimension is high

Cosine Similarity



Very well for
high dimension data

↳ NLP



→ Imputation

Gender	Class	Age	Fruits
M	L		
F	H	-	
M	Med		
F			
M		-	
F		-	

33

29

34

Mean or median

→ Entire data median

→ Female median Age Grouped
Male median Age median

M 45

F -

M 50

✓ F 23

M 55

✓ F 35

M -

$f_1, f_2, f_{j-1}, f_j, f_{j+1}, \dots, f_d$

x^i

	23	
	45	
	—	
	—	

Missing

(Few value are missing)

x

f_j Target

f_1	f_2	...	f_d	f_j
				23
				45
				50
				—
				—
				—

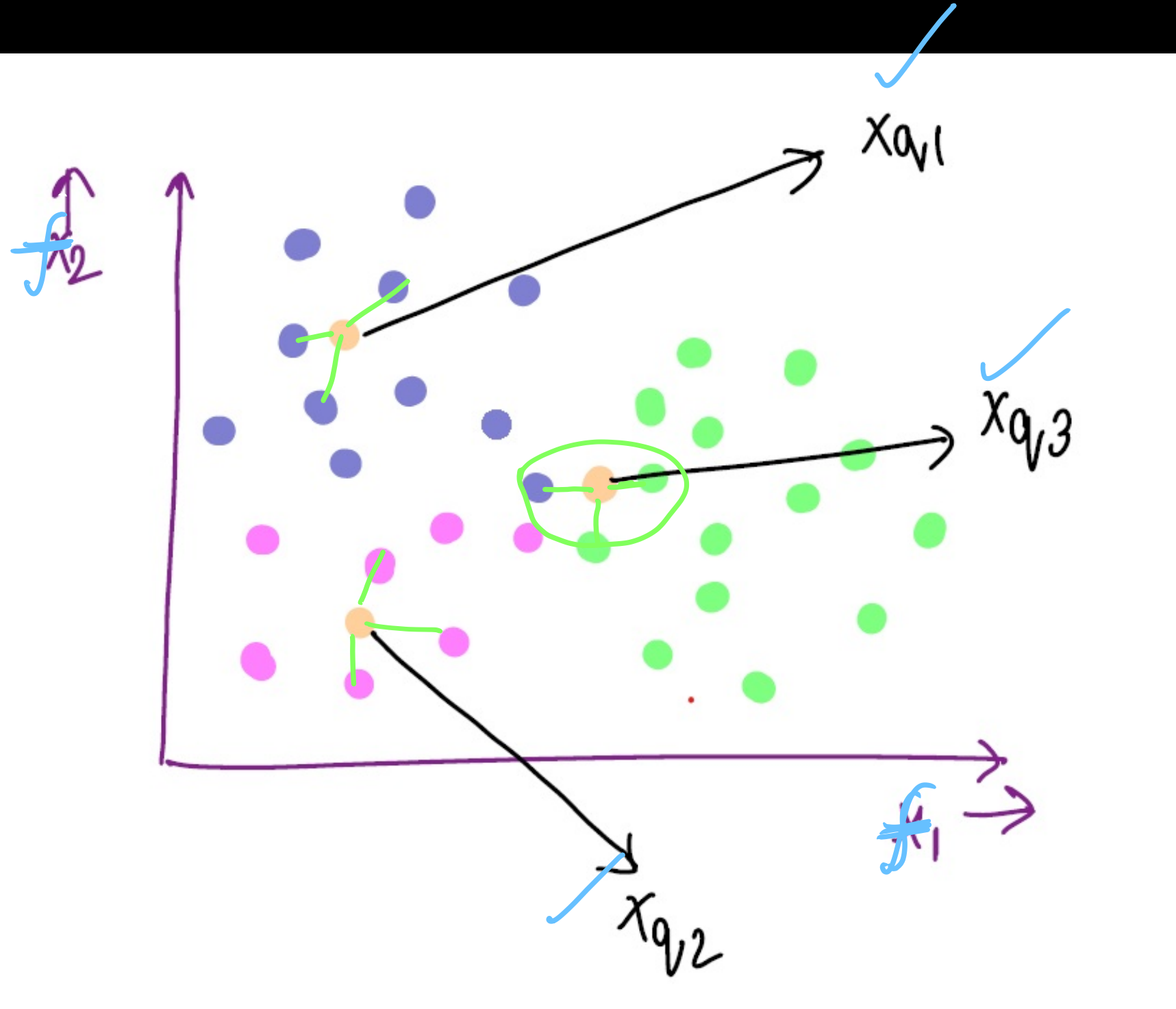
Training

Test

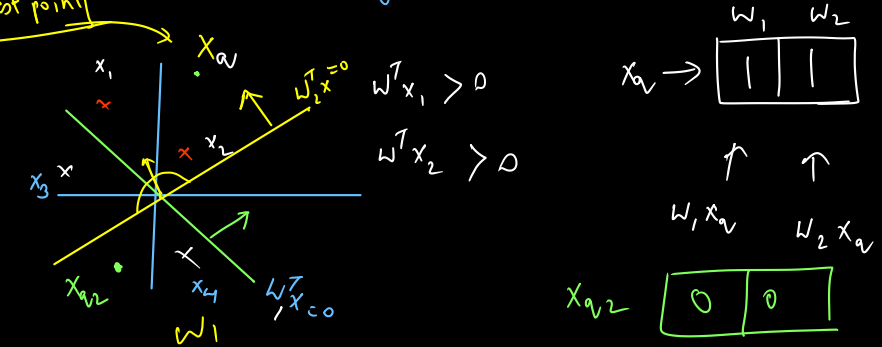
Missing values

x_q

A - blue
B - pink
AB - green



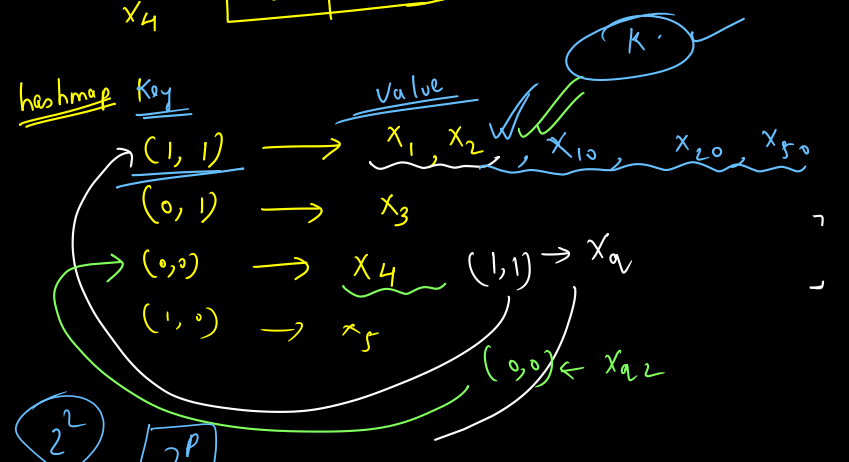
hot point → LSH (Locality Sensitive Hashing) (p.d)



$x_1 \rightarrow w_1^T x_1 > 0, w_2^T x_1 > 0$
 $x_2 \rightarrow w_1^T x_2 > 0, w_2^T x_2 > 0$
 $x_3 \rightarrow w_1^T x_3 < 0, w_2^T x_3 > 0$
 $x_4 \rightarrow w_1^T x_4 < 0, w_2^T x_4 < 0$

	w_1	w_2
x_1	1	1
x_2	1	1
x_3	0	1
x_4	0	0

$w^T x > 0 \rightarrow 1$
 $w^T x < 0 \rightarrow 0$



$O(md + m \log m)$
 $O(n'd)$
 $+ O(p'd)$
 $O(md)$
 \downarrow
 All training data

→ SCAN
 → FAISS

→ Hash table

<u>key</u>	<u>value</u>
{	:
_____	_____
}	