* **Avg word2vec :→**

$$r_{1}: w_1\ w_2\ w_1\ w_3\ w_4\ w_5$$

$n_1$ words

$$V_1 = \frac{1}{n_1}\left[ \underbrace{w_2 V(w_1)}_{d-dim} + \underbrace{w_2 V(w_2)}_{d-dim} + w_2 V(w_1) + \cdots + w_2 V(w_5)\right]$$

* **tf-idf - word2vec :→**

$$r_1: w_1\ w_2\ w_1\ w_3\ w_4\ w_5$$

tfidf :

| $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ | $w_6$ | $w_7$ |
|-------|-------|-------|-------|-------|-------|-------|
| $t_1$ | $t_2$ | $t_3$ | $t_4$ | $t_5$ | $0$ | $0$ |

t stands for tf * idf.

$$tfidf - w_2 v(r_1) = \frac{\sum\limits_{i:words}\left(t_i * w_2 V(w_i)\right)}{\sum\limits_{i:words} t_i}$$

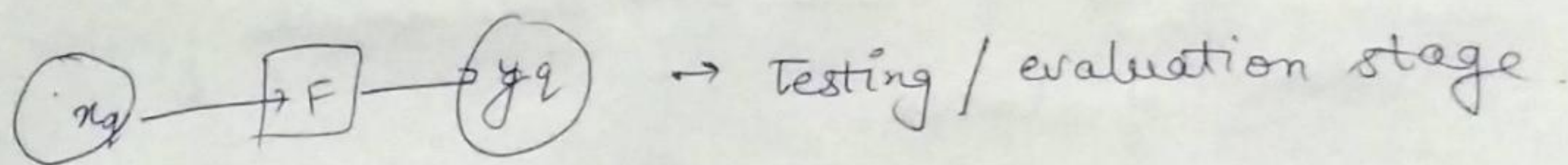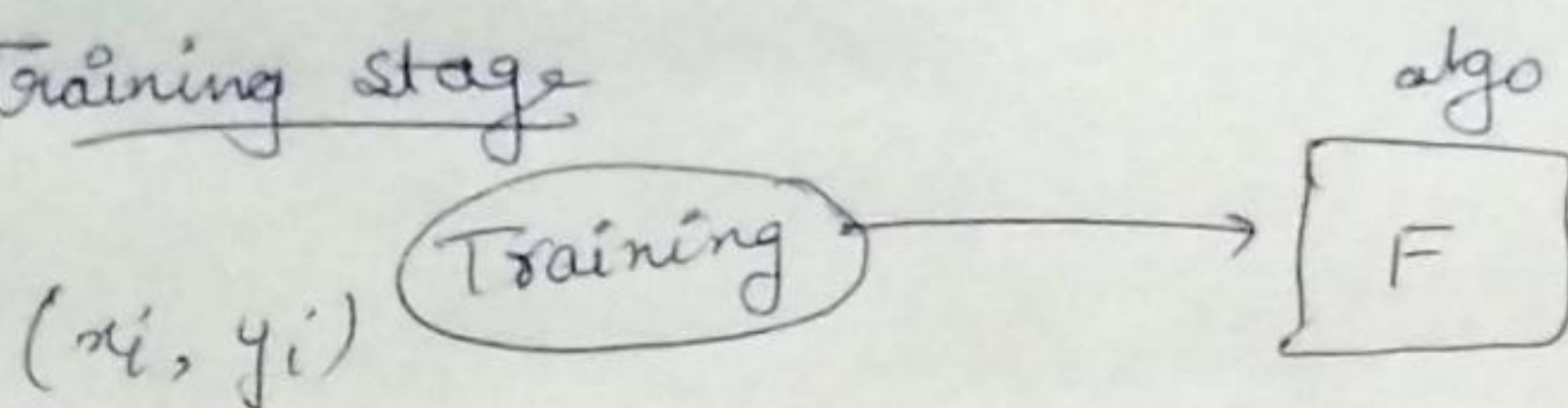* **Classification :→** is to check that given a new review text, determine / predict if the review is positive or negative.

It's like finding a function.

$$\boxed{y = f(x)} \rightarrow \text{central concept.}$$

+ve / -ve

review text

Training stage

(x_i, y_i) (Training) ───────→ algo [ F ]

( n_q ) ─── [→ F ] ─── ( y_q ) → Testing / evaluation stage.

$$D_n = \left\{ (x_i, y_i)_{i=1}^{n} \mid x_i \in R^d, \ y_i \in \{0, 1\} \right\}$$

↓ set.

↓ Such that

↓ ↓
-ve  +ve

Classification algorithm takes these $x_i$ and $y_i$ values as its training data set and then when you provide it any $x_i$ it will apply the function it formed and tell you if the review is positive or not.

* Classification  Vs  Regression :→

When $y_i \in \{0, 1\}$ , when $y_i$ has only two values.

↓ ↓
-ve +ve    then it is called 2 - class

classification / Binary classification.

But in  MNIST dataset,

$$y_i \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

↓

10 class / multiclass classification.

When you've $y_i \in R$

$y_i$ is no more part of a small finite set of classes

then it is called a Regression problem.

So, the only diff in classification & Regression is the value of $y_i$.

i) $y_i \in \{0,1\} \longrightarrow$ or smaller number of values $\longrightarrow$ finite

its a classification problem

ii) $y_i \in R \longrightarrow$ Regression

Starting of KNN

★ K- Nearest Neighbours :→

$x_q \longrightarrow$ k-NN of $x_q$

$\downarrow$

$x_1, x_2, \cdots x_k$

$\downarrow$

$y_1, y_2, \cdots y_k$

$\downarrow$

majority vote

$\downarrow$

$y_q$ :- Ans

* Failure cases of KNN :→

a) when $x_i$ is so far away from all data points then,
its difficult for KNN to solve it.

b) If data is so much concentrated and close. and randomly
be ~~~~ of spread. mixedup data. then,
there is no useful info in it
and KNN fails and most of the algorithm fails.

If KNN works,
apply majority rule.
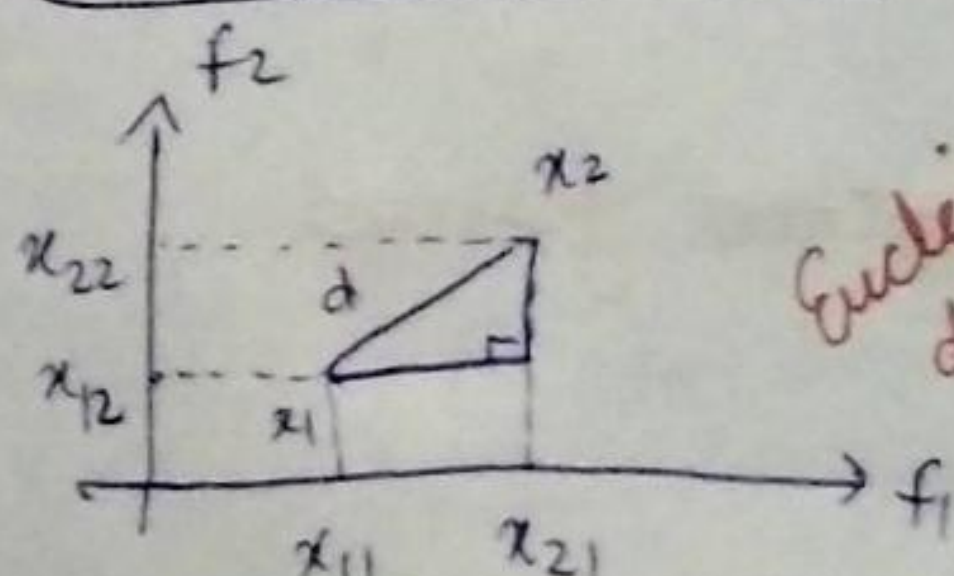
$K = $ odd $\{y_1, y_2, y_3\}$

$+ \quad + \quad + \quad \longrightarrow (+) \rightarrow y_q \rightarrow +ve.$

$+ \quad + \quad - \quad \longrightarrow (+) \rightarrow y_q = +ve$

$\bigotimes$ ii) $K = 4$

$++ -- \longrightarrow (+) \text{ or } (-)$

do not
work in
this case

* Distance Measures :→

$x_1 = (x_{11}, x_{12})$ , $x_2 = (x_{21}, x_{22})$

Euclidean distance

$d = $ len of shortest line from $x_1$ to $x_2$
len

$d = \sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2} = || x_1 - x_2 ||$

for - d- dimensions :→

$$x_i \in R^d , \quad x_2 \in R^d$$

Euclidean distance $= \|x_1 - x_2\|_2 = \left( \sum_{i=1}^{d} \left( x_{1i} - x_{2i} \right)^2 \right)^{1/2}$

↓

$L2$ Norm of a vector

\* **Manhattan dist** :→ $\sum_{i=1}^{d} \underbrace{\left| x_{1i} - x_{2i} \right|}_{\text{absolute value}}$

$\|x_1 - x_2\|_1$ ↗

→ $L_1$ norm of vector $(x_1 - x_1)$

\* **Minkowski dist** :→ generalised form

↓

$L_p$ norms of vector

$$\|x_1 - x_2\|_p = \left( \sum_{i=1}^{d} \left| x_{1i} - x_{2i} \right|^p \right)^{1/p}$$

for $p = 2 \longrightarrow$ minkowski dist $\longrightarrow$ Euclidean dist.

← $\|x_1\|_p = \left( \sum_{i=1}^{d} \left| x_{1i} \right|^p \right)^{1/p}$ , $p \neq 0, p > 0$

$L_p$ norm

⇒ Distances are b/w two points

⇒ Norms are for the corresponding vector formed.

* Hamming dist :→ for boolean vectors.

$x_1, x_2$ ⟶ boolean vector ⟶ Binary Bag of words.

$x_1 = [0, 1, 1, 0, 1, 0, 0 ; \cdots ]$ ⎫
       ↕ ↕           ↕             ⎬ → Hamming dist
$x_2 = [1, 0, 1, 0, 1, 0, 1, \cdots ]$ ⎭    is 3.

Hamming dist $(x_1, x_2)$ = no. of locations / dimensions
                          where binary vectors differ.

Also used for strings.

$x_1 = $ abcade fghik , $x_2 = $ acb ade gf hik

haming dist $(x_1, x_2) = 4$

Hamming dist is useful in case of Gene-code/seq

AGTC

$x_1 = $ A A G T C T C A G . . .
          ↕ ↕       ↕ ↕
$x_2 = $ A G A T C T C G A . . . .

Hamming dist = 4.

* Cosine - similarity & cosine - distance :→

Similarity          distance

↓ dec               ↑ inc            opposite relation.

↑ inc               ↓ dec

$$ \boxed{1 - \cos\text{-}sim\,(x_1, x_2) = \cos\text{-}dist\,(x_1, x_2)} $$

Scanned by TapScanner

let cos-sim $(x_1, x_2) \rightarrow [-1, 1]$

let v. similar     cos-sim $(x_1, x_2) = +1$

v. dissimilar

$$\cos\text{-sim} (x_1, x_2) = -1$$



$d = $ euc-dist.

$$\left\{ \begin{array}{c} \cos\text{-sim} = \cos\theta \\ (x_1, x_2) \end{array} \right\}$$

$\theta$: angle b/w $x_1$ & $x_2$



$$\boxed{\begin{array}{c} \cos \text{ dist } = 1 - \cos\theta \\ (x_1, x_2) \end{array}}$$

cos-sim $(x_1, x_2) = \cos\theta$

cos-sim $(x_1, x_3) = 1$

as $\theta_{x_1, x_3} = 0°$, $\cos 0° = 1$

$$\boxed{\cos\text{-dist}(x_1, x_3) = 1 - 1 = 0}$$
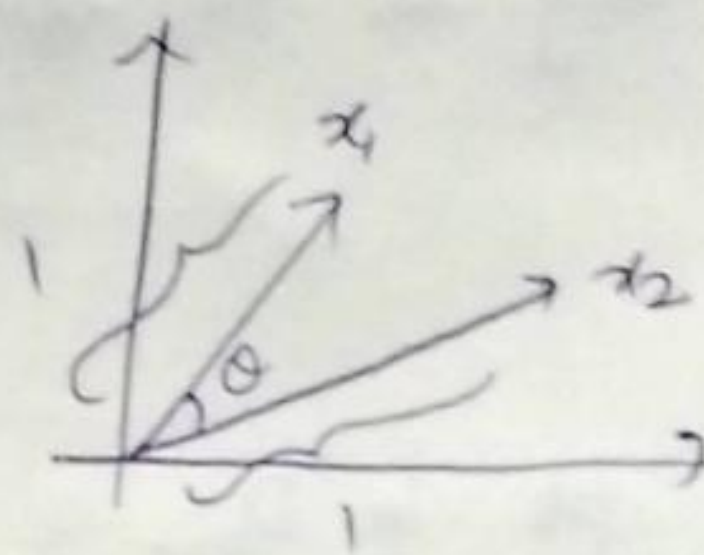
So, even if euclidian distances,

$$\left\{ \begin{array}{c} d_{13} > d_{12} \\ \cos\text{-dist}_{13} < \cos\text{-dist}_{12} \end{array} \right\}$$

angle $(x_1, x_2)$        dist

$0° - 90°$      $\rightarrow$      $0 \rightarrow 1$

$90° - 180°$      $\rightarrow$      $1 \rightarrow 2$

$180° - 270°$      $\rightarrow$      $2 \rightarrow 1$

$270° - 360°$      $\rightarrow$      $1 \rightarrow 0$

$$\cos\theta = \frac{x_1 \cdot x_2}{\|x_1\|_2 \, \|x_2\|_2}$$

cos-sim ↙

$L_2$ norm of $x_1$

① If $x_1$ & $x_2$ are unit vectors

$$\|x_1\|_2 = \|x_2\| = 1$$

$$\boxed{\cos\theta = x_1 \cdot x_2}$$

② Relationship b/w euc-dis & cos-sim ( cos-dist)

if $x_1$ & $x_2$ are unit vectors , where $\theta$ = angle b/w $x_1$ & $x_2$

$$[euc\text{-}dist(x_1, x_2)]^2 = 2\left[1 - \underset{cos\,sim}{\cos(\theta)}\right]$$

$$[euc\text{-}dis(x_1, x_2)]^2 = 2\, cos\text{-}dist(x_1, x_2)$$

\* How to measure how good KNN is?



data points
$$D_n \xrightarrow{\quad} D_{train} \rightarrow n_1$$
$$D_n \searrow D_{test} \rightarrow n_2$$

$$n_1 + n_2 = n$$

How to split $D_n \xrightarrow{R} \overset{70\%}{} D_{train}$
$$D_n \xrightarrow[R]{30\%} D_{test}$$

Randomly

One way is to split it randomly

count = 0;

for each pt in $D_{Test}$:

$$x_q = \text{pt}$$

$$x_q \rightarrow y_1$$

use $D_{Train}$ & k-NN to determine $y_q$

if $y_q == y_{pt}$

count += 1

$\Rightarrow$ count = no. of points for which $D_{train}$ KNN gave a correct class label.

$$Accuracy = \frac{count}{n_2} = \text{no. of pts for which } D_{train}$$
$$+ KNN \text{ gave correct}$$
$$\text{class label.}$$

no of points in $D_{Test}$

$$\boxed{0 \leq Accuracy \leq 1}$$

if Accuracy = 0.91 $\Rightarrow$ 91% of times.

Test - Evaluation time & space complexity 3→

$x_q \longrightarrow y_q$

Input :→ $D_{Train}$ , k, $x_q \in R^d$ ; output : $y_q$

KNN pts = [ ]

O(nd) ⎰ for each $x_i$ in $D_{Train}$ : → n points → can be large enough
⎱ d-dim → can be large enough

O(d) ← Compute $d(x_i, x_q) \longrightarrow d_i$

↓ time complexity          ↓ distance b/w $x_i$ & $x_q$

KNN pts [ ]

O(1) ← Keep the smallest k- distances = $(x_i, y_i, d_i)$

cnt-pos = 0 ; cnt-neg = 0

for each $x_i$ in KNN pts :

small
↑
O(k)
|||
O(1)

     if $y_i$ is +ve

       cnt_pos += 1

     else

       cnt-neg += 1

O(1)

     if cnt_pos > cnt_neg

       return $y_q = 1$ → +ve

     else

       $y_q = 0$ → -ve

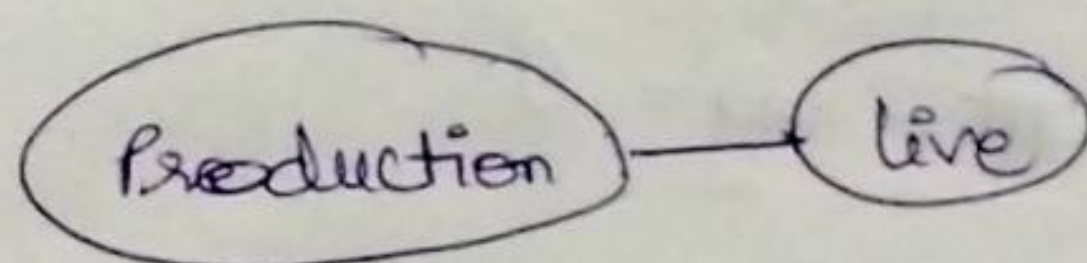Time complexity → $O(nd) + O(1) + O(1)$

$$= \boxed{O(nd)}$$

if $d$ is small

if $d << n$

$$\boxed{O(n)}$$

* Space complexity :→ Space that is needed to evaluate

$$x_q \rightarrow y_q$$

$$\boxed{O(nd)} \rightarrow \text{to store } D_{Train}.$$

* limitations of KNN ( simple implementation)

(Amazon)   fine food Reviews :→

$$\text{Production} - \text{live}$$
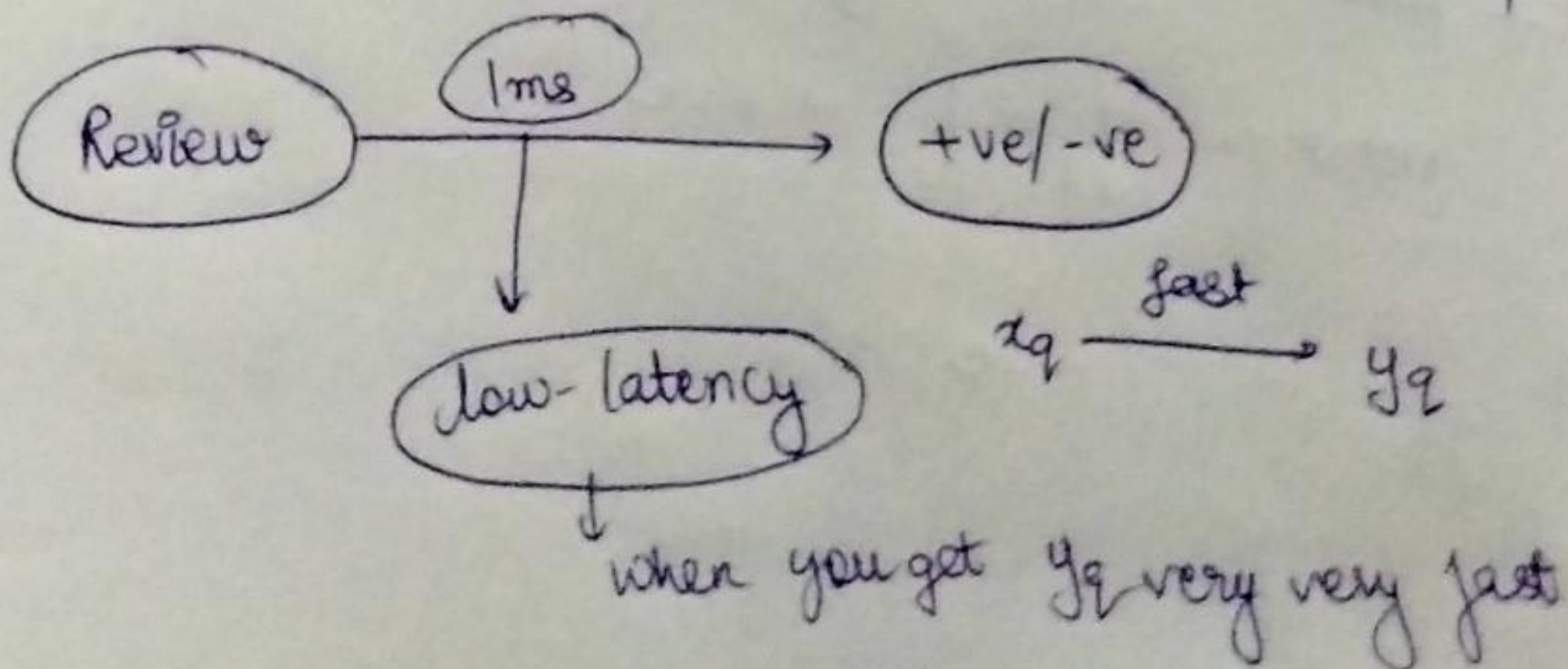
Time :→ $O(nd)$

Space :→ $O(nd)$

$n \approx 364k$
$d \approx 100k$  → DTrain

$36,400M \approx \boxed{36 GB} \rightarrow$ too large space needed.

② Time - complexity → 36 Billions computations.

Review ── 1ms ──→ +ve/-ve

low-latency
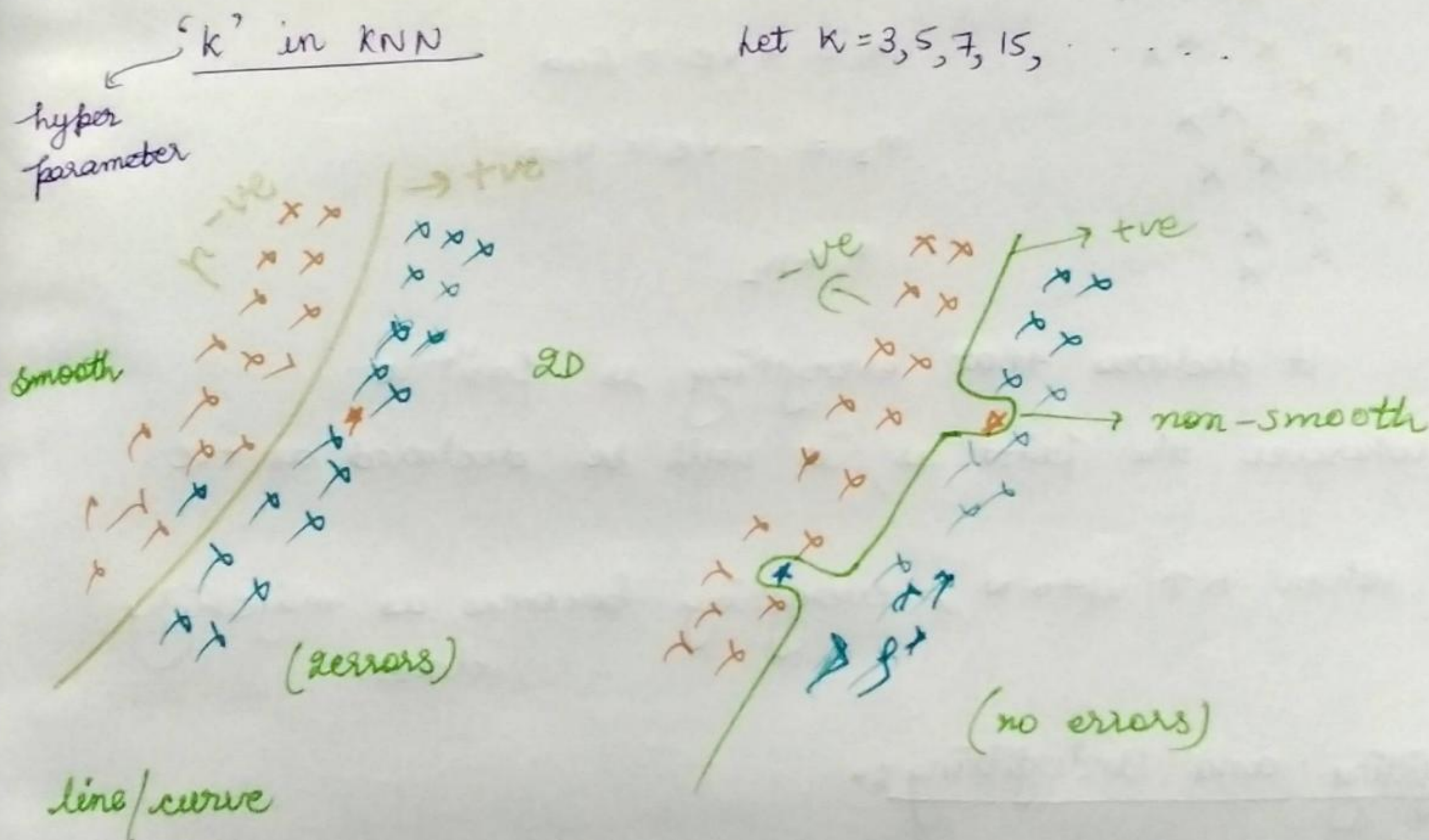
$x_q \xrightarrow{fast} y_q$

$O(nd)$ is too bad
↓
few seconds
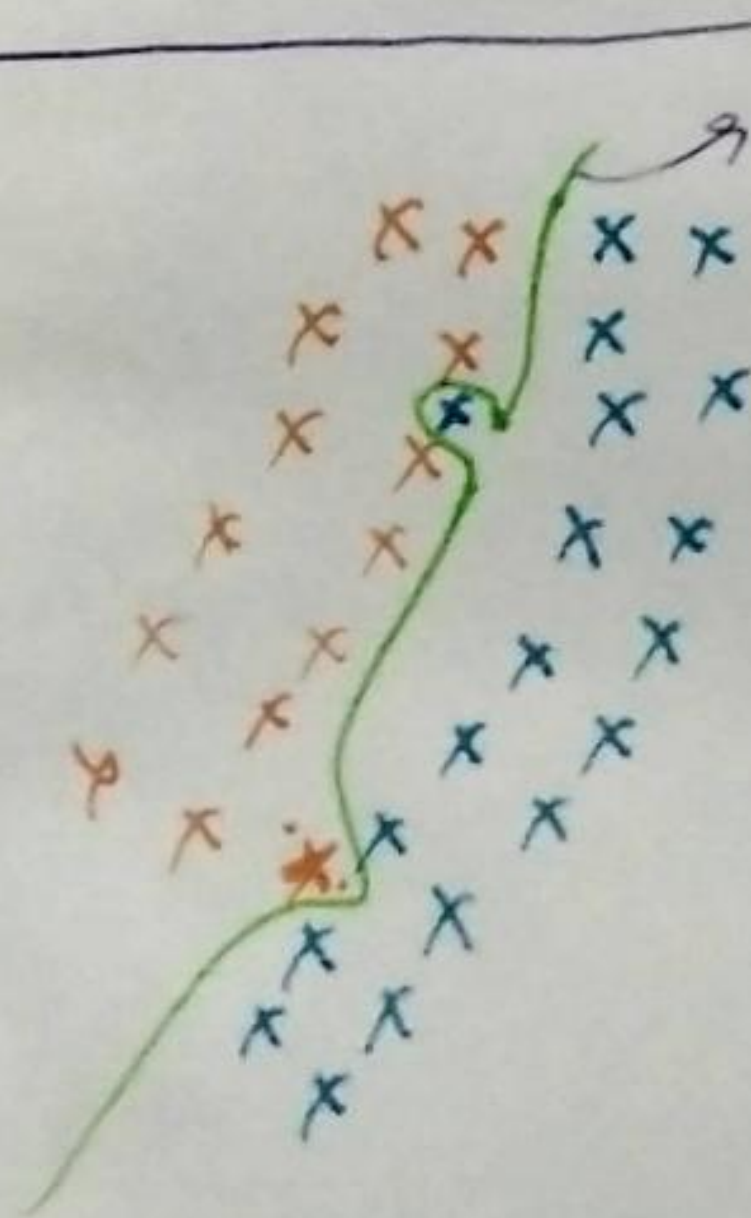
when you get $y_q$ very very fast

So, only reason not to use KNN is its terrible time and space complexity of $O(nd)$.

otherwise its really simple & intuitive & elegant.

* Decision surface for KNN :→

'k' in KNN                    Let K = 3, 5, 7, 15, . . . . .

hyper
parameter

smooth                    2D                              → non-smooth

(2 errors)                                        (no errors)

line/curve

These curves are called decision surfaces.

→ decision surface          ↘ many query points
   when K=1
                            → 1NN, (K=1)

                            → Seperate regions using a
                              curve

                            → smooth
                              curve
                              ↳ Let K=5
                                majority rule.

                            as K↑, smoothness of the
                            curve increases.

So, In KNN, the smoothness of the decision surface increases as K increases.

worst case :→

x x x
x x x x
x
x x x x
x x x x
x x x x
x x x x
x x x x
x x x
x x

1000 - nearest neighbours (NN)

$k = n = 1000$

$n_1 = +ve = 600$

$n_2 = -ve = 400$

$n_1 > n_2$

It declares that everything is positive.

So, wherever the point is, it will be declared as +ve.

So, when K↑ upto N, everything becomes as majority class.