

1) Support Vector Machines ↴

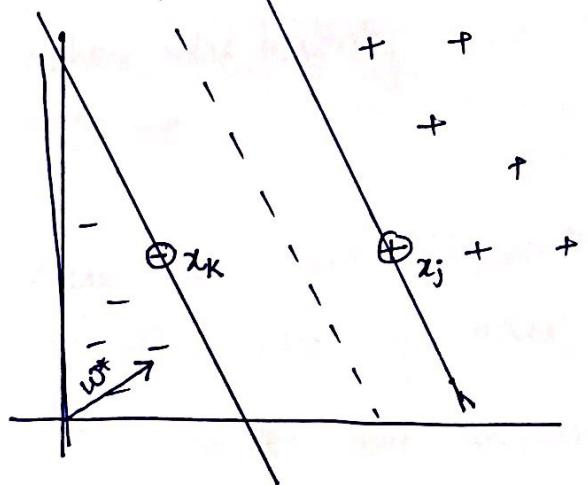
1.1) Hard-margin Linear SVM ↴

Given,

$$\underset{w,b}{\operatorname{argmax}} \quad d = d^+ - d^- = \frac{w^T x_j^+}{\|w\|^2} - \frac{w^T x_k^-}{\|w\|^2}$$

Now given that there are only two support vectors i.e.

$x_j, y_j = +1$ & $x_k, y_k = -1$ with parameters w^* & b^* .



from above, we know that ↴

$$\vec{w}^* \cdot \vec{x}_j + b^* \geq 1 \quad \text{--- (1)}$$

$$\vec{w}^* \cdot \vec{x}_k + b^* \leq -1 \quad \text{--- (2)}$$

Introducing a new variable y_i^* such that

$$y_i^* = \begin{cases} +1 & \text{for +ve samples} \\ -1 & \text{for -ve samples} \end{cases}$$

Multiplying y_i^* with (1) & (2)

from (1), $y_i^* (\vec{w}^* \cdot \vec{x}_j + b^*) \geq 1 \quad \text{such that } y_i^* = +1 \quad \text{--- (3)}$

from (2), $y_i^* (\vec{w}^* \cdot \vec{x}_k + b^*) \leq -1 \quad \text{such that } y_i^* = -1 \quad \text{--- (4)}$

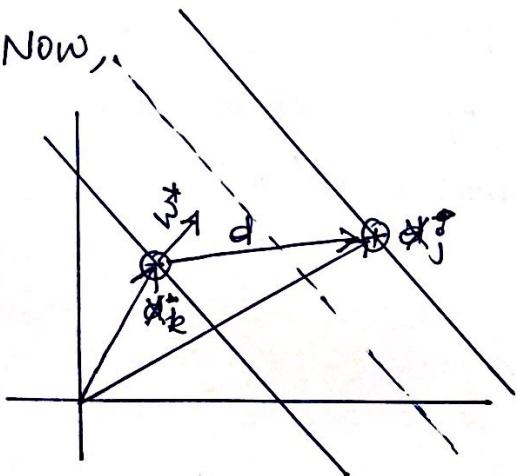
from ③ & ④,

$$y_i (\vec{x}_i \cdot \vec{w}^* + b^*) - 1 \geq 0 \quad \text{such that } y_i = \begin{cases} +1, \\ -1 \end{cases}$$

Also, ③ & ④ can be written as,

$$\left. \begin{array}{l} 1 (\vec{x}_j \cdot \vec{w}^* + b^*) \geq 1 \quad \text{--- (5)} \\ -1 (\vec{x}_j \cdot \vec{w}^* + b^*) \geq 1 \quad \text{--- (6)} \end{array} \right\} \text{Two constraints}$$

Now,



$$d = \frac{(\vec{x}_j^* - \vec{x}_k^*) \cdot \vec{w}^*}{\|\vec{w}^*\|}$$

$$d = (\vec{x}_j^* - \vec{x}_k^*) \cdot \frac{\vec{w}^*}{\|\vec{w}^*\|}$$

from ⑤ & ⑥, substituting the values of \vec{x}_j^* & \vec{x}_k^*

$$d = (1 - \beta^* + 1 + \beta^*) \cdot \frac{1}{\|\vec{w}^*\|}$$

$$\underset{w}{\operatorname{argmax}} \frac{2}{\|\vec{w}^*\|} \Rightarrow \min_{\vec{w}} \|\vec{w}^*\| \Rightarrow \min_{\vec{w}} \frac{1}{2} \|\vec{w}^*\|^2$$

1.2) Soft-Margin SVM

a) In soft-margin SVM, we penalize the incorrectly classified examples by introducing a slack variable i.e. ξ_i . Now, we will minimize the following:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

\$\xi_i\$ ↓
no of mistakes/incorrect classifications
[C = slack penalty]

such that,

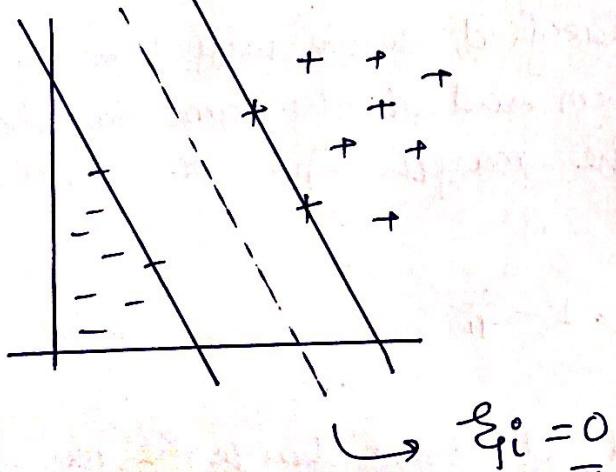
$$y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1 - \xi_i$$

when the training example is correctly classified then

$$\xi_i = 0$$

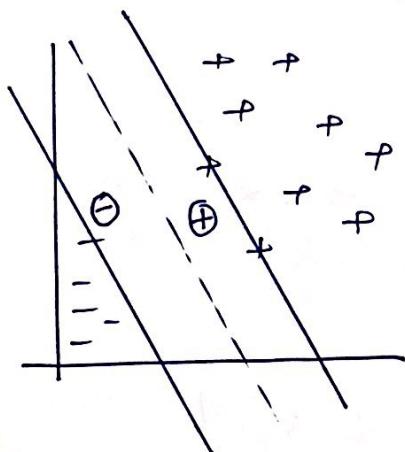
b) there are three different types of examples possible based on the value of slack variable i.e.

i) All examples are correctly classified. In this case, ξ_i will be zero, i.e. penalty will be zero.



Decision boundary will change only if we remove the only training examples that ~~are~~ is on the decision boundary or if that is a support vector.

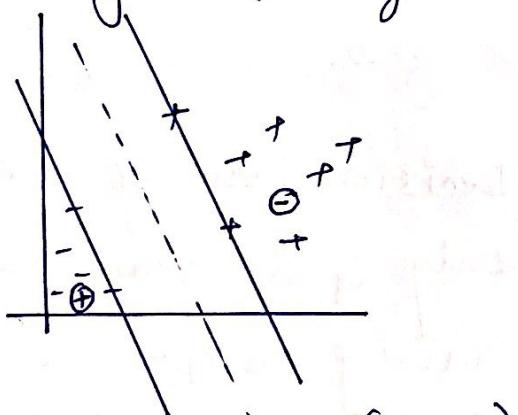
- 2) the second example will be, it is correctly classified but the confidence is low, so in this case the value of ξ_i will be positive but less than 1 i.e. $0 < \xi_i < 1$



here the two datapoints near to the decision boundary are correctly classified but confidence is low, hence they will be penalized, means ξ_i values will be b/w 0 & 1.

Removing these two datapoints will not change the decision boundary/margin.

- 3) the third example will be, when few examples are wrongly classified i.e. $\xi_i > 1$, in this case we will heavily penalize such examples and try to maximize our margin.



here two datapoints are wrongly classified, so we will heavily penalize them and at the same time maximize the margin b/w two classes.

$$(w x + b) y_i \geq 1 - \xi_i$$

Removing such example will not change the margin/decision surface, but removing such example will make the penalty zero and it will be linearly classifiable.

2) Adaboost

2.1) Given, Accuracy = 50%

$$\text{i.e. } \epsilon_t = 0.5$$

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

$$= \frac{1}{2} \ln \left(\frac{0.5}{0.5} \right)$$

$$= \frac{1}{2} \ln(1) = 0$$

Since $\alpha_t = 0$, it means it will not have any say in the final classifier and also, the weight won't change in the next iteration, hence we will drop this ^{weak} classifier and go with some other classifier.

2.2) We are given,

$$\epsilon_t = 0.45 = 0.55$$

here, we can flip the classifier and can get the accuracy equal to 0.55 which is better than the current ϵ_t .

i.e. if now we have,

$$c(x) = \begin{cases} +1 & x > 0 \\ -1 & x \leq 0 \end{cases} \quad \epsilon_t = 0.55$$

then after flipping, we will get,

$$c(x) = \begin{cases} +1 & x \leq 0 \\ -1 & x > 0 \end{cases} \quad \epsilon_t = \underline{\underline{0.45}}$$

i.e. 55% accuracy

Given,

x	0	1	2	3	4	5	6	7	8	9
y	1	1	1	-1	-1	-1	1	1	1	-1

$$C(x) = \begin{cases} +1 & x < 0 \\ -1 & x \geq 0 \end{cases}$$

2.3)

$$\theta = 2.5$$

Initially we have,

$$D_1(0) = D_1(1) = D_1(2) = D_1(3) = D_1(4) = D_1(5) = D_1(6) = D_1(7) = D_1(8) \\ = D_1(9) = \frac{1}{10}$$

for Deriving the weights D_2 ,

$$D_2(i) = \frac{D_1(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

$$\epsilon_t = \frac{1}{\sum_{i=1}^M D_t(i)} \sum_{i=1}^M D_t(i) \delta(h_t(x_i) \neq y_i)$$

$$= 0 + 0 + 0 + 0 + 0 + 0 + \frac{1}{10} \times 1 + \frac{1}{10} \times 1 + \frac{1}{10} \times 1 + 0$$

$$= \frac{3}{10}$$

$$\alpha_t = \frac{1}{2} \ln \left(\frac{7/10}{3/10} \right) = \frac{1}{2} \ln(7/3) = \frac{0.847}{2} = 0.423$$

$$D_2(i) = \frac{D_1(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

$$Z_t = \sum_{i=1}^{10} D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

$$D_2(0) = \frac{\frac{1}{10} \times \exp(-0.423)}{Z_t} = \frac{\frac{1}{10} \times 0.655}{Z_t} = 0.0716$$

$$D_2(1) = \frac{\frac{1}{10} \times \exp(-0.423)}{Z_t} = \frac{\frac{1}{10} \times 0.655}{Z_t} = 0.0716$$

$$D_2(3) = \frac{\frac{1}{10} \times \exp(-0.423)}{Z_t} = \frac{\frac{1}{10} \times 0.655}{Z_t} = 0.0716$$

$$D_2(4) = \frac{\frac{1}{10} \times \exp(-0.423)}{Z_t} = \frac{\frac{1}{10} \times 0.655}{Z_t} = 0.0716$$

$$D_2(5) = \frac{\frac{1}{10} \times \exp(-0.423)}{Z_t} = \frac{\frac{1}{10} \times 0.655}{Z_t} = 0.0716$$

$$D_2(6) = \frac{\frac{1}{10} \times \exp(0.423)}{Z_t} = \frac{\frac{1}{10} \times 1.52}{Z_t} = 0.166$$

$$D_2(7) = \frac{\frac{1}{10} \times \exp(0.423)}{Z_t} = \frac{\frac{1}{10} \times 1.52}{Z_t} = 0.166$$

$$D_2(8) = \frac{\frac{1}{10} \times \exp(0.423)}{Z_t} = \frac{\frac{1}{10} \times 1.52}{Z_t} = 0.166$$

$$D_2(9) = \frac{\frac{1}{10} \times \exp(0.423)}{Z_t} = \frac{\frac{1}{10} \times 0.655}{Z_t} = 0.0716$$

$$Z_t = \frac{1}{10} \times 0.655 \times 9 + \frac{1}{10} \times 1.52 \times 3 = 0.4585 + 0.456 = \underline{\underline{0.9145}}$$

2.4)

$$D_3(i) = \frac{D_2(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

$$\epsilon_t = \frac{1}{\sum_{i=1}^{10} D_t(i)} \sum_{i=1}^{10} D_t(i) \delta(h_t(x_i) \neq y_i)$$

Now for θ , we have 3 possible values, i.e.

$$\theta = 2.5$$

or

$$\theta = 5.5$$

$$\text{or } \theta = 8.5$$

Now for each θ value, we will check for the minimum error and select that θ is decision boundary for our current weak classifier.

$$\underline{\theta = 2.5}$$

$$C(x) = \begin{cases} +1 & x < \theta \\ -1 & x \geq \theta \end{cases}$$

$$\begin{aligned} \epsilon_t &= 0.166 \times 3 \\ &= 0.498 \end{aligned}$$

$$C(x) = \begin{cases} -1 & x < \theta \\ +1 & x \geq \theta \end{cases}$$

$$\begin{aligned} \epsilon_t &= 0.0916 \times 7 \\ &= 0.502 \end{aligned}$$

$$\underline{\theta = 5.5}$$

$$C(x) = \begin{cases} +1 & x < \theta \\ -1 & x \geq \theta \end{cases}$$

$$\begin{aligned} \epsilon_t &= 0.0916 \times 3 + 0.166 \times 3 \\ &= 0.2148 + 0.498 \\ &= 0.7128 \end{aligned}$$

$$C(x) = \begin{cases} +1 & x < \theta \\ -1 & x \geq \theta \end{cases}$$

$$\begin{aligned} \epsilon_t &= 0.0916 \times 4 \\ &= 0.2864 \end{aligned}$$

$$\theta = 0.5$$

$$C(x) = \begin{cases} +1 & x < 0 \\ -1 & x \geq 0 \end{cases}$$

$$C(x) = \begin{cases} +1 & x < 0 \\ +1 & x \geq 0 \end{cases}$$

$$E_t = 0.0716 \times 3$$

$$= 0.214$$

MINIMUM

$$\begin{aligned} E_t &= 0.0716 \times 4 + 0.166 \times 3 \\ &= 0.2864 + 0.498 \\ &= 0.7844 \end{aligned}$$

for all possible 6 error values, we will select the θ & $C(x)$ with minimum error value. i.e

$$\theta = 0.5$$

$$C(x) = \begin{cases} +1 & x < 0 \\ -1 & x \geq 0 \end{cases}$$

} this is our current weak learner for next iteration.

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - 0.214}{0.214} \right)$$

$$= \frac{1}{2} \ln \left(\frac{0.786}{0.214} \right)$$

$$= \frac{1}{2} \ln (3.67) = \frac{1}{2} \times 1.30 = \underline{\underline{0.65}}$$

Now we will calculate $D_3(0) \dots D_3(9)$ using above α_t & $C(x)$.

$$D_3(0) = \frac{D_2(0) \exp(-\alpha + y_0 h_t(0))}{Z_t}$$

$$= \frac{0.0716 \times \exp(-0.65)}{Z_t} = \frac{0.037}{Z_t} = \frac{0.037}{0.817} = 0.045$$

$$D_3(1) = \frac{0.0716 \times \exp(-0.65)}{Z_t} = \frac{0.037}{Z_t} = \frac{0.037}{0.817} = 0.045$$

$$D_3(2) = \frac{0.0716 \times \exp(-0.65)}{Z_t} = \frac{0.037}{Z_t} = \frac{0.037}{0.817} = 0.045$$

$$D_3(3) = \frac{0.0716 \times \exp(0.65)}{Z_t} = \frac{0.137}{Z_t} = \frac{0.137}{0.817} = 0.167$$

$$D_3(4) = \frac{0.0716 \times \exp(0.65)}{Z_t} = \frac{0.137}{Z_t} = \frac{0.137}{0.817} = 0.167$$

$$D_3(5) = \frac{0.0716 \times \exp(0.65)}{Z_t} = \frac{0.137}{Z_t} = \frac{0.137}{0.817} = 0.167$$

$$D_3(6) = \frac{0.166 \times \exp(-0.65)}{Z_t} = \frac{0.086}{Z_t} = \frac{0.086}{0.817} = 0.105$$

$$D_3(7) = \frac{0.166 \times \exp(-0.65)}{Z_t} = \frac{0.086}{Z_t} = \frac{0.086}{0.817} = 0.105$$

$$D_3(8) = \frac{0.166 \times \exp(-0.65)}{Z_t} = \frac{0.086}{Z_t} = \frac{0.086}{0.817} = 0.105$$

$$D_3(9) = \frac{0.0716 \times \exp(-0.65)}{Z_t} = \frac{0.037}{Z_t} = \frac{0.037}{0.817} = 0.045$$

$$\begin{aligned} Z_t &= 0.037 \times 4 + 0.086 \times 3 + 0.137 \times 3 \\ &= 0.148 + 0.258 + 0.411 = 0.817 \end{aligned}$$

for D_4 , again we have 3 possible values for θ i.e. $2.5, 5.54$
 8.5 .

Now, we will calculate error at each θ & will choose the θ with min error.

$\theta = 2.5$

$$C(x) = \begin{cases} +1 & x < 0 \\ -1 & x \geq 0 \end{cases}$$

$$E_t = 0.105 \times 3 = 0.315$$

$$C(x) = \begin{cases} -1 & x < 0 \\ +1 & x \geq 0 \end{cases}$$

$$E_t = 0.045 \times 4 + 0.167 \times 3 = 0.681$$

$\theta = 5.5$

$$C(x) = \begin{cases} +1 & x < 0 \\ -1 & x \geq 0 \end{cases}$$

$$\begin{aligned} E_t &= 0.167 \times 3 + 0.105 \times 3 \\ &= 0.504 + 0.315 = 0.819 \end{aligned}$$

$$C(x) = \begin{cases} -1 & x < 0 \\ +1 & x \geq 0 \end{cases}$$

$$\begin{aligned} E_t &= 0.045 \times 4 \\ &= \boxed{0.18} \text{ MINIMUM} \end{aligned}$$

$\theta = 8.5$

$$C(x) = \begin{cases} +1 & x < 0 \\ -1 & x \geq 0 \end{cases}$$

$$E_t = 0.167 \times 3 = 0.501$$

$$C(x) = \begin{cases} -1 & x < 0 \\ +1 & x \geq 0 \end{cases}$$

$$\begin{aligned} E_t &= 0.045 \times 4 + 0.105 \times 3 \\ &= 0.18 + 0.315 = 0.49 \end{aligned}$$

Out of all 6 possible error, we will select the minimum i.e

$\theta = 5.5$

$$C(x) = \begin{cases} -1 & x < 0 \\ +1 & x \geq 0 \end{cases}$$

$$X_t = \frac{1}{2} \ln \left(\frac{1-0.18}{0.18} \right) = \frac{1}{2} \times 1.504 = \underline{\underline{0.752}}$$

$$D_4(i) = \frac{D_3(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

$$D_4(0) = \frac{0.045 \exp(0.75)}{Z_t} = \frac{0.045 \times 2.11}{Z_t} = \frac{0.095}{Z_t} = 0.124$$

$$D_4(1) = \frac{0.045 \exp(0.75)}{Z_t} = \frac{0.045 \times 2.11}{Z_t} = \frac{0.095}{Z_t} = 0.124$$

$$D_4(2) = \frac{0.045 \exp(0.75)}{Z_t} = \frac{0.045 \times 2.11}{Z_t} = \frac{0.095}{Z_t} = 0.124$$

$$D_4(3) = \frac{0.167 \exp(-0.75)}{Z_t} = \frac{0.167 \times 0.472}{Z_t} = \frac{0.078}{Z_t} = 0.102$$

$$D_4(4) = \frac{0.167 \exp(-0.75)}{Z_t} = \frac{0.167 \times 0.472}{Z_t} = \frac{0.078}{Z_t} = 0.102$$

$$D_4(5) = \frac{0.167 \exp(-0.75)}{Z_t} = \frac{0.167 \times 0.472}{Z_t} = \frac{0.078}{Z_t} = 0.102$$

$$D_4(6) = \frac{0.015 \exp(-0.75)}{Z_t} = \frac{0.015 \times 0.472}{Z_t} = \frac{0.049}{Z_t} = 0.064$$

$$D_4(7) = \frac{0.015 \exp(-0.75)}{Z_t} = \frac{0.015 \times 0.472}{Z_t} = \frac{0.049}{Z_t} = 0.064$$

$$D_4(8) = \frac{0.015 \exp(-0.75)}{Z_t} = \frac{0.015 \times 0.472}{Z_t} = \frac{0.049}{Z_t} = 0.064$$

$$D_4(9) = \frac{0.045 \exp(0.75)}{Z_t} = \frac{0.045 \times 2.11}{Z_t} = \frac{0.095}{Z_t} = 0.124$$

$$\begin{aligned} Z_t &= 0.095 \times 4 + 0.078 \times 3 + 0.049 \times 3 \\ &= 0.38 + 0.234 + 0.148 = 0.762 \end{aligned}$$

Q.5)

After four iterations, we have

	$\theta = 2.5$ $+1 \leftarrow 1 \rightarrow -1$			$\theta = 5.5$ $-1 \leftarrow 1 \rightarrow +1$			$\theta = 8.5$ $+1 \leftarrow 1 \rightarrow -1$		
X	0	1	2	3	4	5	6	7	8 9
Y	1	1	1	-1	-1	-1	1	1	1 -1

from above we can see that our data has been classified correctly, hence classification error is zero.

So our final classifier will have those 3 weak classifiers (decision tree stumps) weighted with their α_t values i.e.

$$H(x) = \text{Sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

Since our classifier have reach 100% classification accuracy, we can stop the iteration and take weighted sum of above three classifiers as our final classifier.

3) K-Nearest Neighbour classifier ↴

3.1) A Lazy Classifier :-

a) When a new training example becomes available, we have to retrain the models from scratch for Naïve Bayes classifier and Support Vector Machines. For both the algorithm we have to set up our hyper-parameters based on the new training example as these new incoming training example can change the decision boundary/surface over-time. So to get this new boundary/surface, we have to retrain our model from scratch and learn new hyper-parameters.

Whereas, for KNN, since it is a instance based classifier, and there is no model or hyper-parameters are involved, new training example will not make any difference. But when new test example comes, we have to calculate the distance from all examples including new training example.

b) KNN will need the most computation to infer the class label for this new test example as it will have to calculate the distance of new test ^{example} from each of the ~~each~~ training example to find the K-nearest neighbours and use them to label the new example.

Time complexity $\rightarrow O(n) + O(n) \xrightarrow{\text{for distance}} \xrightarrow{\text{for sorting (using Quicksort)}} O(n)$

while SVM & Naïve Bayes will take $O(1)$ time to classify new example.

3.2) Implement KNN

a) Step 1 → pre-process the data. Since we have used just first 10% of the data for training set and test set.

Step 2 → for each of the test example, calculate its distance from each of the training example, and get the top k-nearest neighbours for that particular example. We have used the Euclidean distance as our distance metric.

e.g. A B
 (x_1, y_1) (x_2, y_2)

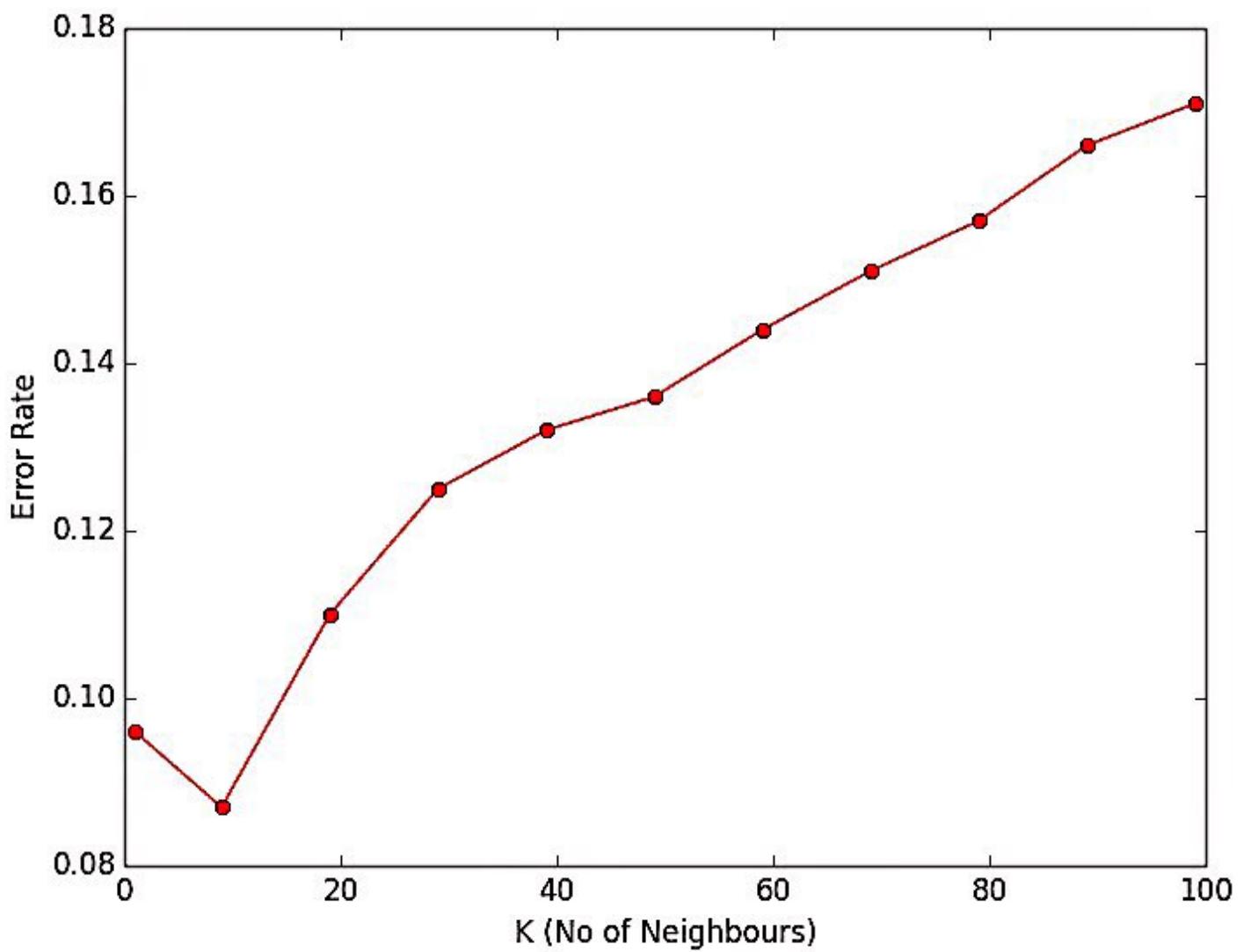
$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

In general,

$$d = \sqrt{\sum_{i,j}^n (x_i - x_j)^2}$$

Step 3 → Now that we have the top-k neighbours, we will predict the label of our test example by taking the vote of each of the top-k neighbour's label. Highest voted label will be assigned to the new test example.

Step 4 → Accuracy can be calculated based on the actual label vs the predicted label for particular test example.



```
For K = 1
Test Error: 0.096
No of Incorrect Labels: 96
-----
For K = 9
[Test Error: 0.086
No of Incorrect Labels: 86
[-----
For K = 19
Test Error: 0.11
No of Incorrect Labels: 110
-----
[For K = 29
Test Error: 0.123
No of Incorrect Labels: 123
-----
For K = 39
Test Error: 0.132
No of Incorrect Labels: 132
-----
For K = 49
Test Error: 0.136
No of Incorrect Labels: 136
-----
For K = 59
Test Error: 0.143
No of Incorrect Labels: 143
-----
For K = 69
Test Error: 0.151
No of Incorrect Labels: 151
-----
For K = 79
Test Error: 0.156
No of Incorrect Labels: 156
-----
For K = 89
Test Error: 0.166
No of Incorrect Labels: 166
-----
For K = 99
Test Error: 0.171
No of Incorrect Labels: 171
```