

# Nearest Neighbor Classification

The Nearest-Neighbor Rule

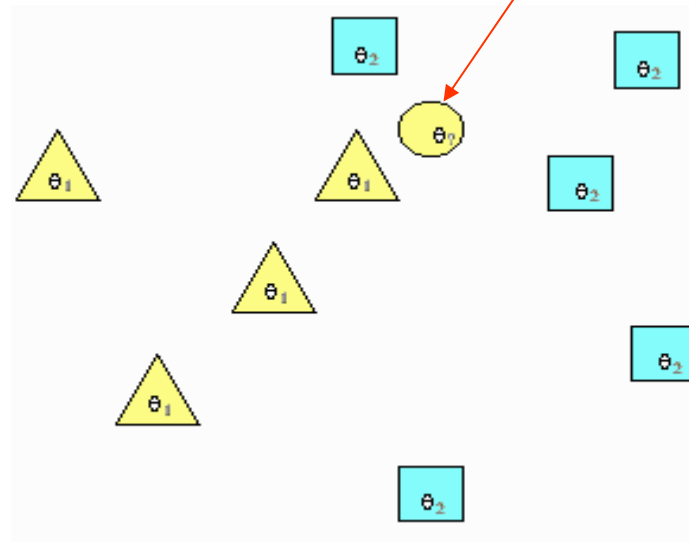
Error Bounds

k-Nearest Neighbor Rule

Computational Considerations

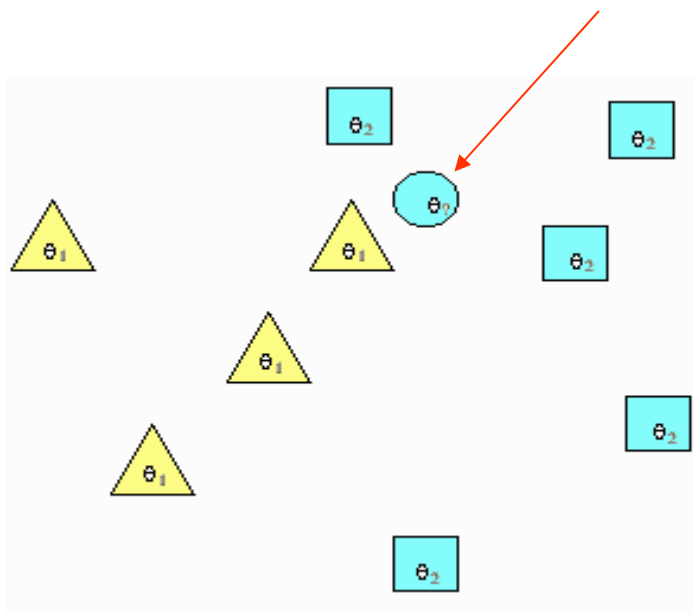
# Example of Nearest Neighbor Rule

- Two class problem: yellow triangles and blue squares. Circle represents the unknown sample  $x$  and as its nearest neighbor comes from class  $\theta_1$ , it is labeled as class  $\theta_1$ .



# Example of k-NN rule with $k = 3$

- There are two classes: yellow triangles and blue squares. The circle represents the unknown sample  $x$  and as two of its nearest neighbors come from class  $\theta_2$ , it is labeled class  $\theta_2$ .

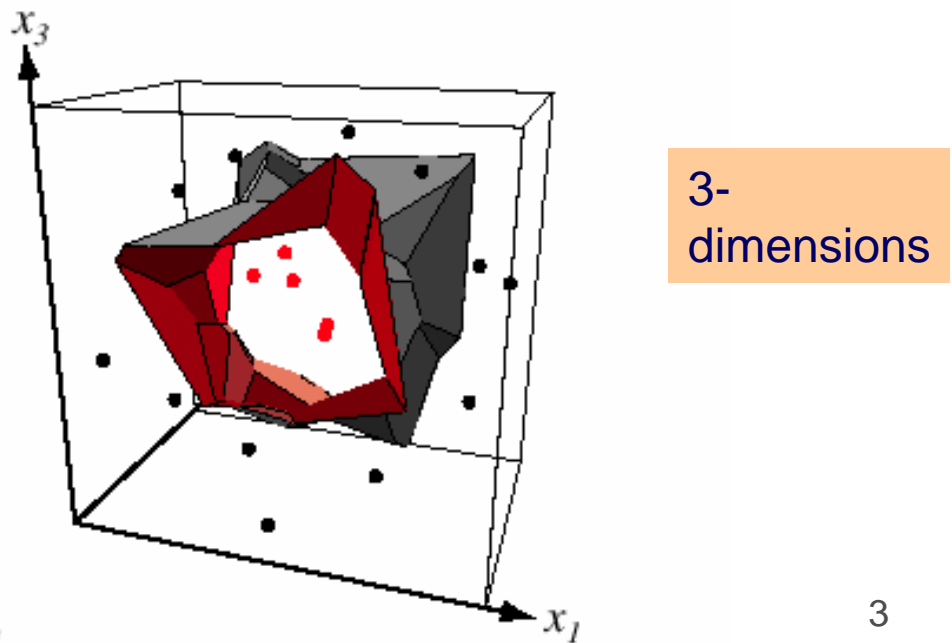
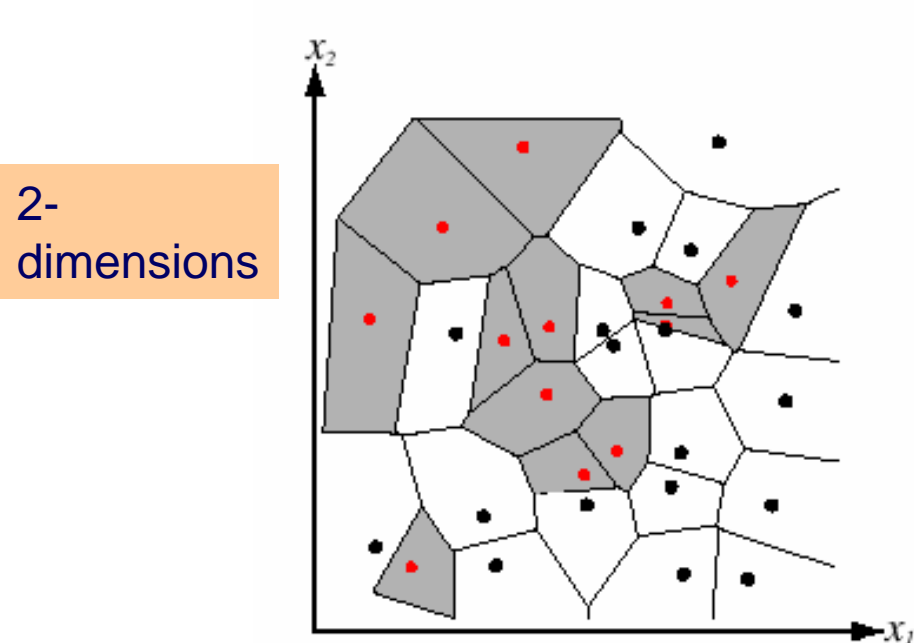


The number  $k$  should be:

- 1) large to minimize probability of misclassifying  $x$ .
- 2) small (with respect to no of samples) so that points are close enough to  $x$  to give an accurate estimate of the true class of  $x$ .

# Nearest Neighbor and Voronoi Tesselation

- N-n classifier effectively partitions the feature space into cells consisting of all points closer to a given training point  $x'$  than to any other training points.
- All points in such a cell are thus labeled by the category of the training point – *Voronoi tessellation* of the space



# Nearest–Neighbor Rule Probability of Error

Let  $D^n = \{x_1, x_2, \dots, x_n\}$  be a set of  $n$  labeled prototypes

- Let  $x' \in D^n$  be the nearest prototype to a test point  $x$
- The nearest-neighbor rule for classifying  $x$  is
  - to assign it the label associated with  $x'$
- Nearest-neighbor rule is a sub-optimal procedure
  - Does not yield the Bayes error rate
  - Yet it is never worse than twice the Bayes error rate

# Why does Nearest–Neighbor rule work well?

- Label  $\theta'$  associated with nearest neighbor is a random variable
- Probability that  $\theta' = \omega_i$  is the *a posteriori* probability  $P(\omega_i / x')$
- As  $n \rightarrow \infty$ , it is always possible to find  $x'$  sufficiently close so that:

$$P(\omega_i / x') \cong P(\omega_i / x)$$

- Because this is exactly the probability that nature will be in state  $\omega_i$  the nearest neighbor rule is effectively matching probabilities with nature

# Bayesian Probability of Error

- If we define  $\omega_m(\mathbf{x})$  by

$$P(\omega_m | \mathbf{x}) = \max_i P(\omega_i | \mathbf{x}),$$

then the Bayes decision rule always selects  $\omega_m$ .

From this the Bayesian condition probability of error is

$$P^*(e | x) = 1 - P(\omega_m | x)$$

# Bayesian Probability of Error

If we let  $P^*(e/x)$  be the minimum possible value of  $P(e/x)$ , and  $P^*$  be the minimum possible value of  $P(e)$ , then by averaging over the a priori distribution of  $x$  we get

$$P^* = \int P^*(e | x) p(x) dx = \int (1 - P(\omega_m | x)) p(x) dx$$



# Evaluation of Nearest Neighbor Error

- If  $P_n(e)$  is the  $n$  - sample error rate, and if

$$P = \lim_{n \rightarrow \infty} P_n(e),$$

Then we want to show that

$$P^* \leq P \leq P^* \left( 2 - \frac{c}{c-1} P^* \right).$$

# Nearest-Neighbor Probability of Error

## The Random Variables

Begin by looking at all the random variables in the construction of an  $x, x_n, \theta, \theta_n$  system.

We denote  $\theta$  as the true class of  $x$  and  $\theta_n$  as the labeled class of  $x_n$ , where  $x_n$  is the nearest neighbor of  $x$ .

It is clear that  $x$  and its  $\theta$  are random input parameters to the problem. Note that the underlying statistics of the labeled space are random too. Thus the  $x_n, \theta_n$  pair are also unknown and thus random inputs.

The probability of  $x$  having true class  $\theta$  and that of  $x_n$  being labeled  $\theta_n$  are independent.

Thus we have

$$P(\theta, \theta_n \mid x, x_n) = P(\theta \mid x)P(\theta_n \mid x_n)$$

# Expressing the Probability of Error

Now if we use the nearest-neighbor decision rule, we commit an error whenever  $\theta \neq \theta'_n$ . Thus, the conditional probability of error  $P_n(e|\mathbf{x}, \mathbf{x}'_n)$  is given by

$$\begin{aligned} P_n(e|\mathbf{x}, \mathbf{x}'_n) &= 1 - \sum_{i=1}^c P(\theta = \omega_i, \theta'_n = \omega_i | \mathbf{x}, \mathbf{x}'_n) \\ &= 1 - \sum_{i=1}^c P(\omega_i | \mathbf{x}) P(\omega_i | \mathbf{x}'_n). \end{aligned} \quad (43)$$

# Convergence of Probability of Error

- Notice that as  $n$  approaches infinity the space of labeled items will become increasingly filled. Thus the nearest neighbor of  $x$  will become  $x_n$  with probability 1.
- So we can say that:

$$\lim_{n \rightarrow \infty} P(e | x, x_n) = \lim_{n \rightarrow \infty} P(e | x, x) = \lim_{n \rightarrow \infty} P(e | x)$$

# Final Expression for Nearest-Neighbor Probability of Error

Therefore, provided that we can exchange some limits and integrals, the asymptotic nearest-neighbor error rate is given by

$$\begin{aligned} P &= \lim_{n \rightarrow \infty} P_n(e) \\ &= \lim_{n \rightarrow \infty} \int P_n(e|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ &= \int \left[ 1 - \sum_{i=1}^c P^2(\omega_i|\mathbf{x}) \right] p(\mathbf{x}) d\mathbf{x}. \end{aligned} \tag{45}$$

# Bounds on the Conditional Probability of Error

While Eq. 45 presents an exact result, it is more illuminating to obtain bounds on  $P$  in terms of the Bayes rate  $P^*$ . An obvious lower bound on  $P$  is  $P^*$  itself. Furthermore,

it can be shown that for any  $P^*$  there is a set of conditional and prior probabilities for which the bound is achieved, so in this sense it is a tight lower bound.

The problem of establishing a tight upper bound is more interesting. The basis for hoping for a low upper bound comes from observing that if the Bayes rate is low,  $P(\omega_i|\mathbf{x})$  is near 1.0 for some  $i$ , say  $i = m$ . Thus the integrand in Eq. 45 is approximately  $1 - P^2(\omega_m|\mathbf{x}) \simeq 2(1 - P(\omega_m|\mathbf{x}))$ , and since

$$P^*(e|\mathbf{x}) = 1 - P(\omega_m|\mathbf{x}), \quad (46)$$

integration over  $\mathbf{x}$  might yield about twice the Bayes rate, which is still low and acceptable for some applications. To obtain an exact upper bound, we must find out how large the nearest-neighbor error rate  $P$  can become for a given Bayes rate  $P^*$ . Thus, Eq. 45 leads us to ask how small  $\sum_{i=1}^c P^2(\omega_i|\mathbf{x})$  can be for a given  $P(\omega_m|\mathbf{x})$ .

# Nearest Neighbor Error Bound Derivation

First we write

$$\sum_{i=1}^c P^2(\omega_i|\mathbf{x}) = P^2(\omega_m|\mathbf{x}) + \sum_{i \neq m} P^2(\omega_i|\mathbf{x}), \quad (47)$$

and then seek to bound this sum by minimizing the second term subject to the following constraints:

- $P(\omega_i|\mathbf{x}) \geq 0$
- $\sum_{i \neq m} P(\omega_i|\mathbf{x}) = 1 - P(\omega_m|\mathbf{x}) = P^*(e|\mathbf{x})$ .

With a little thought we see that  $\sum_{i=1}^c P^2(\omega_i|\mathbf{x})$  is minimized if all of the *a posteriori* probabilities except the *m*th are equal. The second constraint yields

$$P(\omega_i|\mathbf{x}) = \begin{cases} \frac{P^*(e|\mathbf{x})}{c-1} & i \neq m \\ 1 - P^*(e|\mathbf{x}) & i = m. \end{cases} \quad (48)$$

Thus we have the inequalities

$$\sum_{i=1}^c P^2(\omega_i|\mathbf{x}) \geq (1 - P^*(e|\mathbf{x}))^2 + \frac{P^{*2}(e|\mathbf{x})}{c-1} \quad (49)$$

and

$$1 - \sum_{i=1}^c P^2(\omega_i|\mathbf{x}) \leq 2P^*(e|\mathbf{x}) - \frac{c}{c-1} P^{*2}(e|\mathbf{x}). \quad (50)$$



## Error Bound Conclusion

This immediately shows that  $P \leq 2P^*$ , since we can substitute this result in Eq. 45 and merely drop the second term. However, a tighter bound can be obtained by observing that the variance is:

$$\begin{aligned}\text{Var}[P^*(e|\mathbf{x})] &= \int [P^*(e|\mathbf{x}) - P^*]^2 p(\mathbf{x}) d\mathbf{x} \\ &= \int P^{*2}(e|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} - P^{*2} \geq 0,\end{aligned}$$

so that

$$\int P^{*2}(e|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \geq P^{*2}, \quad (51)$$

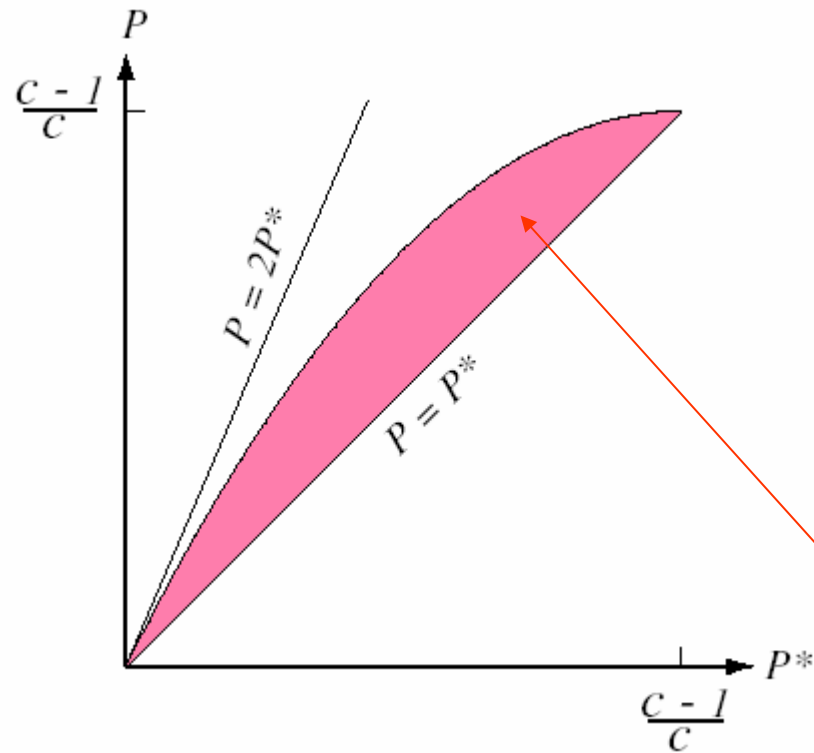
with equality holding if and only if the variance of  $P^*(e|\mathbf{x})$  is zero. Using this result and substituting Eq. 50 into Eq. 45, we obtain the desired bounds on the nearest-neighbor error  $P$  in the case of an infinite number of samples:

$$P^* \leq P \leq P^* \left( 2 - \frac{c}{c-1} P^* \right). \quad (52)$$

Error bounds are tight in that for any  $P^*$  there exist Conditional and prior distributions for which the Bounds are achieved.



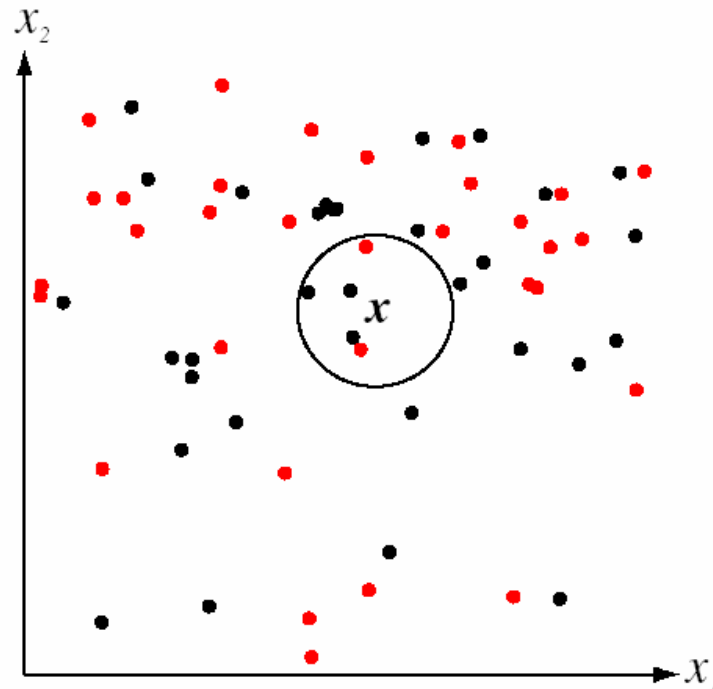
# Bounds on nearest neighbor error rate in c-category problem



Possible  
Asymptotic  
Error rates

# The $k$ – Nearest-Neighbor Rule

- Classify  $x$  by assigning it the label most frequently represented among the  $k$  nearest samples and use a voting scheme



# Analysis of $k$ – Nearest-Neighbor Rule

- Select  $w_m$  if a majority of the  $k$  nearest neighbors are labeled  $w_m$ , an event of probability

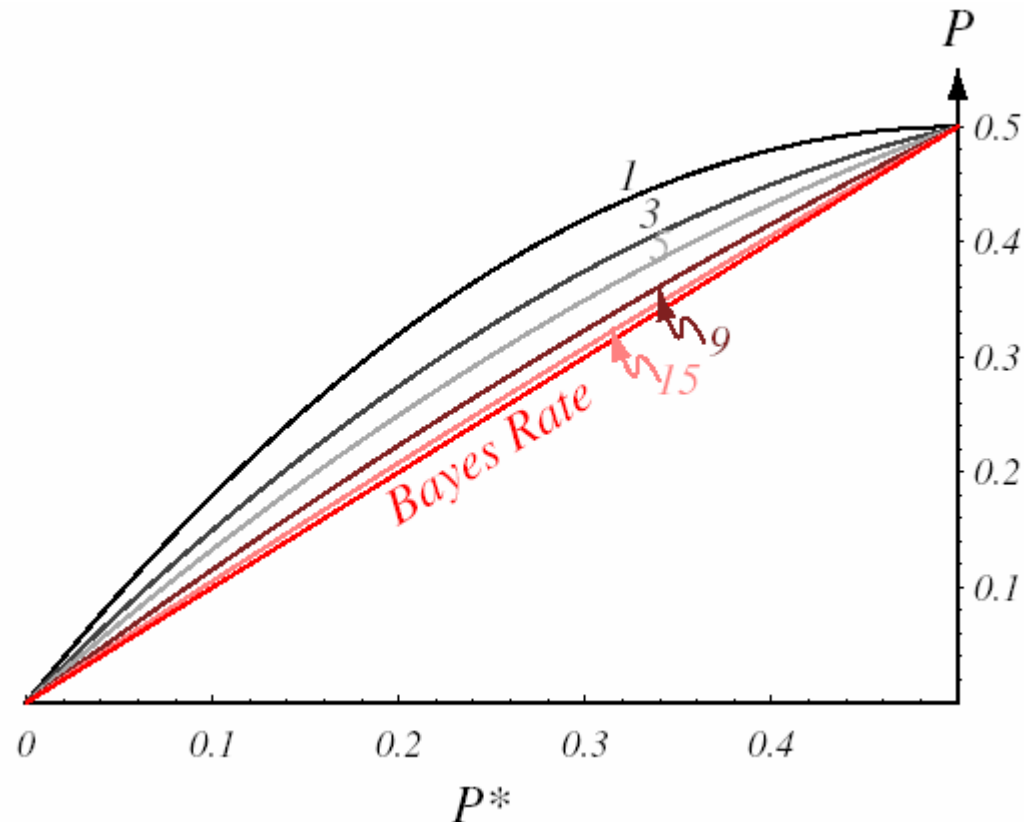
$$\sum_{i=(k+1)/2}^k \binom{k}{i} P(\omega_m|\mathbf{x})^i [1 - P(\omega_m|\mathbf{x})]^{k-i}.$$

- It can be shown that if  $k$  is odd, the large-sample two-class error rate for the  $k$ -nearest-neighbor rule is bounded above by the function  $C_k(P^*)$ , where  $C_k(P^*)$  is defined to be the smallest concave function of  $P^*$  greater than

$$\sum_{i=0}^{(k-1)/2} \binom{k}{i} [(P^*)^{i+1} (1 - P^*)^{k-i} + (P^*)^{k-i} (1 - P^*)^{i+1}].$$

# Bounds on Error Rate of $k$ -Nearest Neighbor Rule

Bound is  
 $C_k(P^*)$

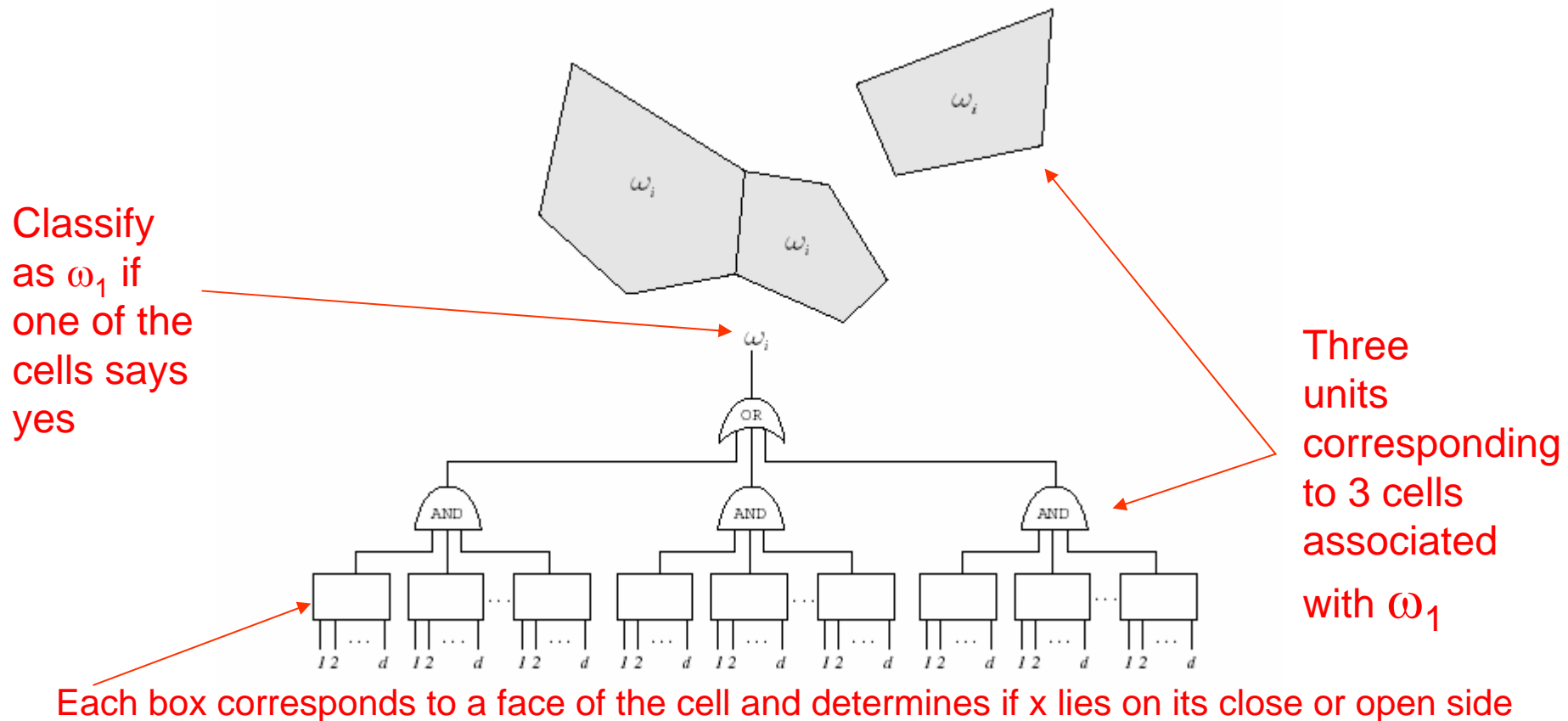


As  $k$  gets larger the error rate equals the Bayes rate  
 $k$  should be a small fraction of the total number of samples

# Computational Complexity of $k$ -Nearest-Neighbor Rule

- Each Distance Calculation is  $O(d)$
- Finding single nearest neighbor is  $O(n)$
- Finding  $k$  nearest neighbors involves sorting; thus  $O(dn^2)$
- Methods for speed-up:
  - Parallelism
  - Partial Distance
  - Pre-structuring
  - Editing, pruning or condensing

# Parallel Implementation of $k$ -Nearest-Neighbor Rule



# Partial Distance Method of n-n speedup

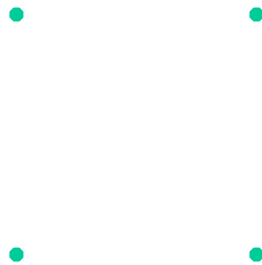
- The *partial distance* based on  $r$  selected dimensions is

$$D_r(\mathbf{a}, \mathbf{b}) = \left( \sum_{k=1}^r (a_k - b_k)^2 \right)^{1/2}$$

- Terminate a distance calculation once its partial distance is greater than the full  $r = d$  Euclidean distance to the current closest prototype

# Search Tree Method of nn speedup

- Create a search tree where prototypes are selectively linked
- Consider only the prototypes linked to entry point



- Points in neighboring region may actually be closer
- Tradeoff of accuracy versus speed



# Editing Method of nn speedup

- Eliminate Prototypes that are surrounded by training points of the same category

## ■ Algorithm 3. (Nearest-Neighbor Editing)

```
1 begin initialize  $j \leftarrow 0$ ,  $\mathcal{D} \leftarrow$  data set,  $n \leftarrow$  # prototypes
2   construct the full Voronoi diagram of  $\mathcal{D}$ 
3   do  $j \leftarrow j + 1$ ; for each prototype  $\mathbf{x}'_j$ 
4     find the Voronoi neighbors of  $\mathbf{x}'_j$ 
5     if any neighbor is not from the same class as  $\mathbf{x}'_j$ , then mark  $\mathbf{x}'_j$ 
6   until  $j = n$ 
7   discard all points that are not marked
8   construct the Voronoi diagram of the remaining (marked) prototypes
9 end
```

- Complexity is  $O(d^3 n^{d/2} \ln n)$