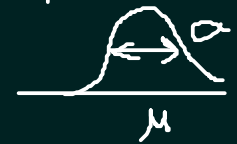


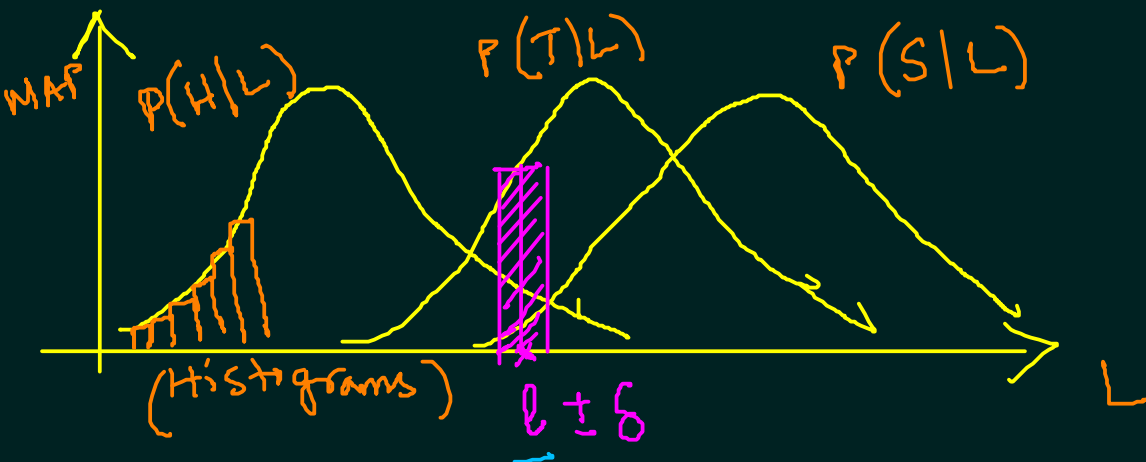
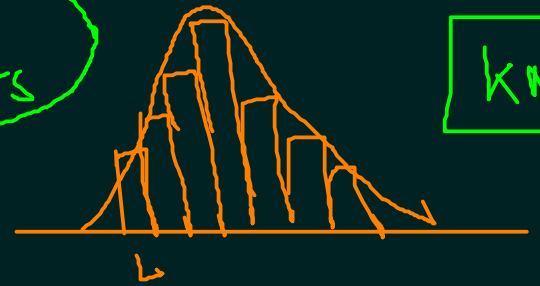
Parametric Learning: MAP \rightarrow MLE, Prior \rightarrow parameters θ



Non-parametric Learning:

k-Nearest Neighbours

KNN



$l = 2ft$
 $\delta = 0.1$
 $1.9ft \leftrightarrow 2.1ft$

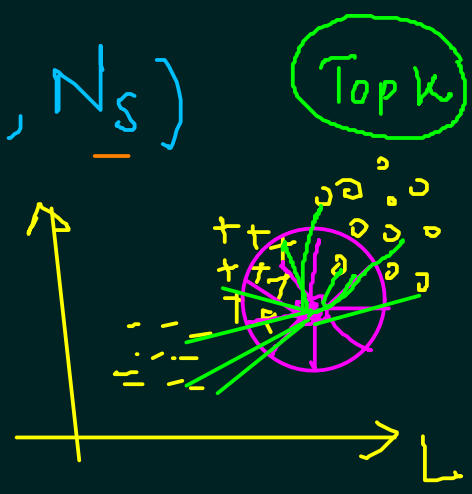
$N = N_H + N_T + N_S$

Choose class for "x" { vote and wins } \rightarrow $f_H = \frac{N_H}{N}$, $f_T = \frac{N_T}{N}$, $f_S = \frac{N_S}{N}$

majority (N_H, N_T, N_S)

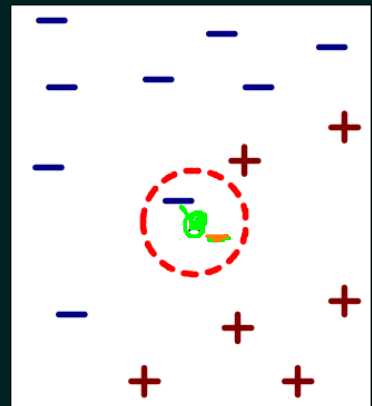
[Approximate MAP \rightarrow (relative count of elements)]

class \leftarrow [majority (N_H, N_T, N_S) \leftarrow $\frac{K_H + K_T + K_S}{k} = k$ rad = δ

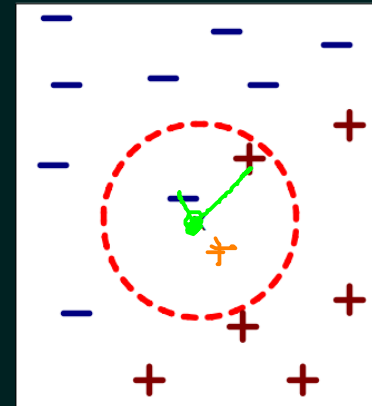


distance between points

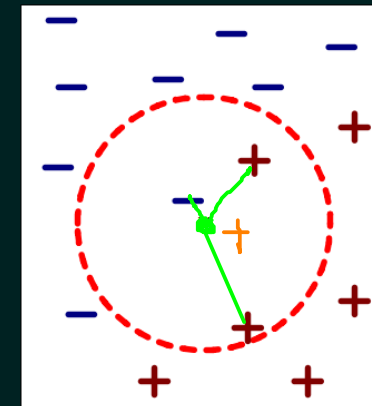
Nearest Neighbour Problem



(a) 1-nearest neighbor $k=1$



(b) 2-nearest neighbor $k=2$



(c) 3-nearest neighbor $k=3$

Voronoi Diagram:

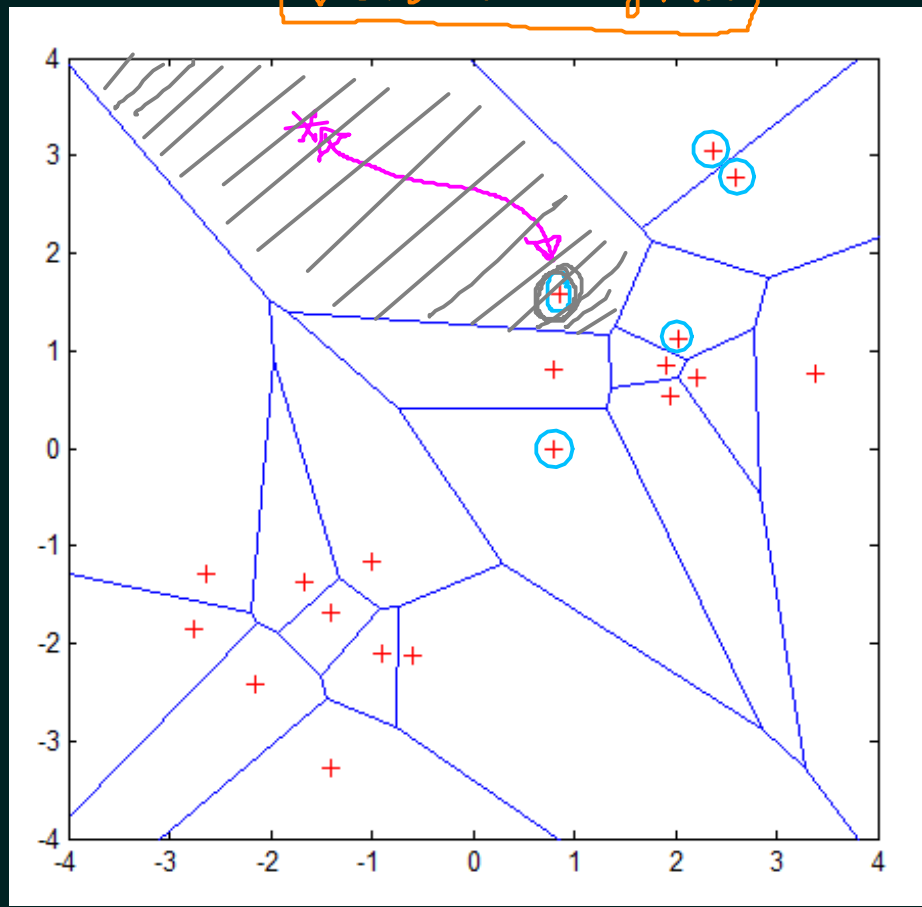
- Tessellation Cell

- Any unknown point
↳ which cell it belongs

- Algorithm wise
Straight fwd.

↳ Train → ✓
↳ Test → Simple } ⊕ to your friend

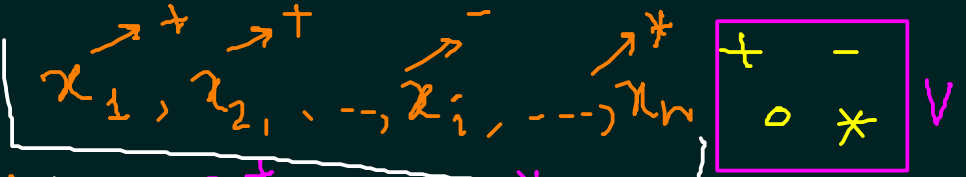
Voronoi Diagram



Distance weighted k-NN

$$\hat{f}(q) = \underset{v \in V}{\operatorname{argmax}} \left(\sum_{j=1}^k \delta(v, f(x_{i_j})) \right)$$

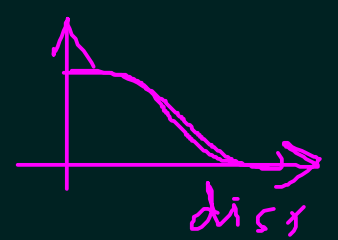
KNN



$$\begin{cases} v=0 \rightarrow \text{count} = \\ v=* \rightarrow \text{count} = \\ v=+ \rightarrow \text{count} = \end{cases}$$

$f(x_{i_1}) = +$ $\delta(a, b) = 1, a=b$
 $f(x_{i_2}) = -$ $= 0, a \neq b$
 $f(x_{i_3}) = 0$ $f(x_{i_4}) = 0$

$$\hat{f}(q) = \underset{v \in V}{\operatorname{argmax}} \sum_{j=1}^k \frac{1}{d(q, x_{i_j})^2} \delta(v, f(x_{i_j}))$$



$$\hat{f}(q) = \underset{v \in V}{\operatorname{argmax}} \sum_{j=1}^k w_j \delta(v, f(x_{i_j}))$$

Continuous values

$$\hat{f}(q) = \frac{\sum_{j=1}^k w_j f(x_{i_j})}{\sum_{j=1}^k w_j}$$

Weighted Avg

What is a good 'k' = ?

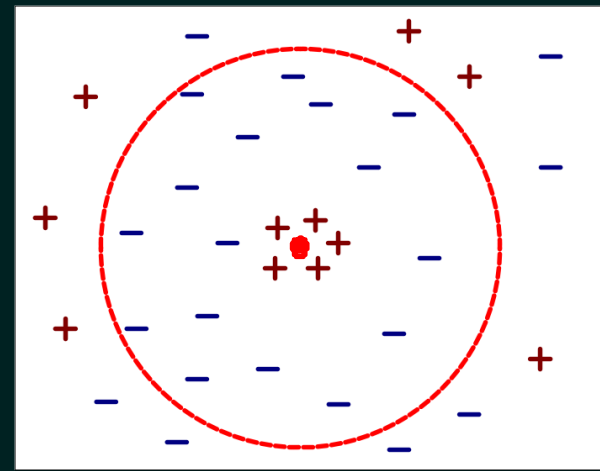
Small $k \rightarrow$ too local
 \rightarrow noise \times

$k = N$ (all TE)

too global \rightarrow prior (MAP), MLE (data likelihood)

Typical : $k = \sqrt{N}$, $N = \#TE$
 $k = \frac{N}{10}$ } Practical est.

(OPEN PROBLEM) \checkmark



Minkowsky:

$$D(x, y) = \left(\sum_{i=1}^m |x_i - y_i|^r \right)^{1/r}$$

Euclidean:

$$D(x, y) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2}$$

Manhattan / city-block:

$$D(x, y) = \sum_{i=1}^m |x_i - y_i|$$

Camberra:

$$D(x, y) = \sum_{i=1}^m \frac{|x_i - y_i|}{|x_i + y_i|}$$

Chebychev:

$$D(x, y) = \max_{i=1}^m |x_i - y_i|$$

Quadratic:

$$D(x, y) = (x - y)^T Q (x - y) = \sum_{j=1}^m \left(\sum_{i=1}^m (x_i - y_i) q_{ji} \right) (x_j - y_j)$$

Q is a problem-specific positive definite $m \times m$ weight matrix

Mahalanobis:

$$D(x, y) = [\det V]^{1/m} (x - y)^T V^{-1} (x - y)$$

V is the covariance matrix of $A_1..A_m$, and A_j is the vector of values for attribute j occurring in the training set instances $1..n$.

Correlation:

$$D(x, y) = \frac{\sum_{i=1}^m (x_i - \bar{x}_i)(y_i - \bar{y}_i)}{\sqrt{\sum_{i=1}^m (x_i - \bar{x}_i)^2 \sum_{i=1}^m (y_i - \bar{y}_i)^2}}$$

$\bar{x}_i = \bar{y}_i$ and is the average value for attribute i occurring in the training set.

Chi-square:

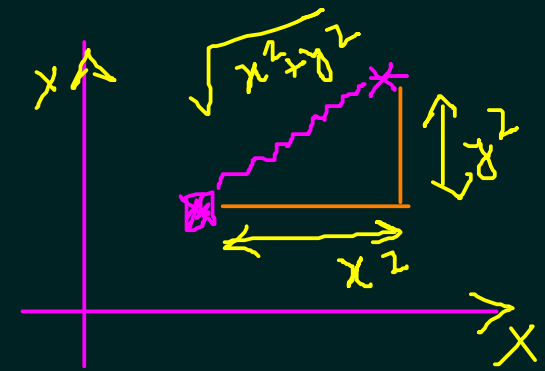
$$D(x, y) = \sum_{i=1}^m \frac{1}{sum_i} \left(\frac{x_i}{size_x} - \frac{y_i}{size_y} \right)^2$$

sum_i is the sum of all values for attribute i occurring in the training set, and $size_x$ is the sum of all values in the vector x .

Kendall's Rank Correlation:

$$D(x, y) = 1 - \frac{2}{n(n-1)} \sum_{i=1}^m \sum_{j=1}^{i-1} \text{sign}(x_i - x_j) \text{sign}(y_i - y_j)$$

$\text{sign}(x) = -1, 0$ or 1 if $x < 0, x = 0,$ or $x > 0,$ respectively.



- Euclidean
- Manhattan
- Minkowsky

↳ Correlation

Distance Metrics

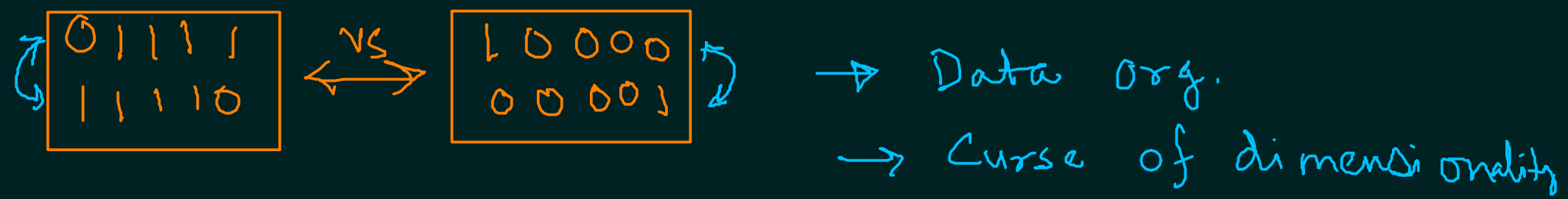
Figure 1. Equations of selected distance functions. (x and y are vectors of m attribute values).

Issue 1

$$q = \begin{matrix} a_1 & a_2 & a_3 & a_4 & a_5 \\ 0 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 0 \end{matrix} \quad \left. \vphantom{\begin{matrix} a_1 & a_2 & a_3 & a_4 & a_5 \\ 0 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 0 \end{matrix}} \right\} \text{Euclidean}$$

$$q = \begin{matrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{matrix}$$

$\Rightarrow d = \sqrt{2}$ $d = \sqrt{2}$



Issue 2

$a_1 = \text{age } [0 - 100]$ $a_3 = \text{bp}$
 $a_2 = \text{blood sugar } [\quad]$ $a_4 = [0.01 - 0.09]$

\rightarrow Normalize

$$X_i = (x_{ij})$$

$$Z_{ij} = \left(\frac{x_{ij} - \mu_i}{\sigma_i} \right)$$

$\mu_i = \text{Mean}(X_i)$
 $\sigma_i = \text{SD}(X_i)$

$N = 20$

$k = \lfloor \sqrt{N} \rfloor = 4$ (odd)



$d(x_1, x_2) = \frac{3}{4}$

Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
whale	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
turtle	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
dolphin	yes	no	yes	no	mammals
eagle	no	yes	no	yes	non-mammals

Give Birth	Can Fly	Live in Water	Have Legs	Class
yes	no	yes	no	?

$d(x_1, q) = 1/2$
 $d(x_2, q) = 1/2$
 $d(x_3, q) = 1/4$
 $d(x_4, q) = 0$

Top k

2-class problem
 $\rightarrow k = \text{odd}$

Algo. Indicates NON-MAMMEL