

Übung 9 Bioinformatik

1. Training Data

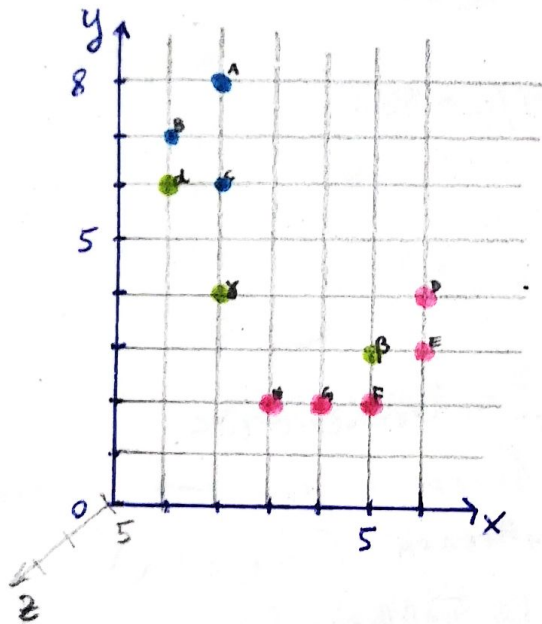
Tumorgewebe:

	A	B	C
x	2	1	2
y	8	7	6
z	5	5	5

gesundes Gewebe

	D	E	F	G	H
x	6	6	5	4	3
y	4	3	2	2	2
z	5	5	5	5	5

dreidimensionale Werte/Punkte sind gegeben, dabei handelt es sich eigentlich um 2 dimensionale Werte, da der z-Wert nicht variiert \rightarrow alle liegen in Ebene \rightarrow 2D



Vorhersagen für

	d	β	γ
x	1	5	2
y	6	3	4
z	5	5	5

k nearest neighbours mit geringstem Abstand finden \rightarrow die Eigenschaft / Merkmale das überwiegt gilt auch für d, β , γ

$K=1$

(d) nearest neighbours: B oder C (eindeutig) jeweils Abstand = 1
beide Tumor (Tu) $\rightarrow d = Tu$

(β) n.n: E oder F (eindeutig erkennbar) $d_{\beta, E} = d_{\beta, F} = 1$
~~beide~~ ~~beide~~ ~~beide~~ beide gesundes Gewebe (gG) $\rightarrow \beta = gG$

(γ) n.n: C oder H $d_{\gamma, C} = \sqrt{0^2 + 2^2} = 2$
 $d_{\gamma, H} = \sqrt{1^2 + 2^2} = \sqrt{5}$
 $\rightarrow d_{\gamma, C} < d_{\gamma, H} \rightarrow \gamma = Tu$

Übung 9 Bioinformatik Seite 2

k=2

① $n.n = B \& C$ (siehe k=1) $\rightarrow d = Tu$

② $n.n = E \& F$ (") $\rightarrow \beta = gG$

③ $n.n = C \& H$ $d_{g,C} = \sqrt{4} = 2$ $\rightarrow \gamma = ?$
 $d_{g,H} = \sqrt{5}$

k=3

① $n.n = A, B, C$ $\rightarrow d = Tu$

② $n.n = E, F, D/G$
 (alles noch sehr eindeutig) $\rightarrow \beta = gG$

③ $n.n = C \& H + B$ oder G (?) $d_{g,G} = \sqrt{2^2 + 2^2} = \sqrt{8}$
 $d_{g,B} = \sqrt{3^2 + 1^2} = \sqrt{10}$
 $d_{g,G} < d_{g,B} \rightarrow \gamma = gG$

k=4

① $n.n = A, B, C, H$ $\rightarrow d = Tu$

② $n.n = D, E, F, G$ $\rightarrow \beta = gG$

③ $n.n = C, H, B, G$ $\rightarrow \gamma = ?$

k=5

① $n.n = A, B, C, H \& G$ oder D

$d_{d,G} = \sqrt{3^2 + 4^2} = \sqrt{25}$

$d_{d,D} = \sqrt{5^2 + 2^2} = \sqrt{29}$

$d_{d,G} < d_{d,D} \rightarrow d = Tu$

② $n.n = D, E, F, G, H$ $\beta \rightarrow gG$

③ $n.n = C, H, B, G \& A/D$ oder F

$d_{g,D} = 4$

$d_{g,F} = \sqrt{3^2 + 2^2} = \sqrt{13}$

$d_{g,D} > d_{g,F} \rightarrow \gamma = gG$

k=6

① $n.n = A, B, C + 3g.G$ $\rightarrow d = ?$

② $n.n = D, E, F, G, H +$ irgendein Tu $\rightarrow \beta = gG$

③ $n.n = C, H, B, G, F, A/D$ $\rightarrow \gamma = gG$ oder ? (je nachdem, ob A/D)

$$k=7$$

$$\textcircled{\alpha} \quad n.n = A, B, C, \dots + 4 g_G \rightarrow \alpha = g_G$$

$$\textcircled{\beta} \quad n.n = D, E, F, G, H, C, B \rightarrow \beta = g_G$$

$$\textcircled{\gamma} \quad n.n = A, B, C, D, E, G, H \rightarrow \gamma = g_G$$

Bewertung: - ^{bei} ~~ab~~ $k=6$ kann es nur noch zu einem "uneabgeschlossen" zwischen g_G & T_u kommen oder zum Überwiegen von g_G , da es nur 3 Trainingsdaten zu T_u gab

- ab $k=7$ ~~es kommt~~ ^{kommt man} ausschließlich zu dem Ergebnis, dass die ^{Mehrheit der} nearest neighbours gesund ist, da diese zahlenmäßig überwiegen

\Rightarrow es kann bei zu hohen k s zu Fehlinterpretationen kommen (wenn diese über der ~~doppelten~~ zweifachen Größe der kleinsten Trainingsdatengruppe liegen)

\Rightarrow ~~es~~ k muss / sollte passend zur Anzahl der Trainingsdatenpunkte gewählt werden

\Rightarrow zudem ungerade, da es bei geraden k s zu nicht aussagekräftigen Ergebnissen kommen kann
(vgl. $\gamma(k=2)$, $\gamma(k=4)$, $\gamma(k=6)$, $\alpha(k=6)$)

2.

$$p(x) = \frac{1}{1 + e^{(-\beta_0 + \beta_1 x)}} = \frac{1}{1 + e^{-z}} \quad \begin{matrix} \beta_0 = -5 \\ \beta_1 = 1 \end{matrix}$$

wurden mit
maximum
Likelihood
Methode berechnet

$$p(2) = \frac{1}{1 + e^{(5+2)}} = 9,11 \cdot 10^{-4}$$

$$p(6) = \frac{1}{1 + e^{(5+6)}} = 1,67 \cdot 10^{-5}$$

$$p(10) = \frac{1}{1 + e^{(5+10)}} = 3,06 \cdot 10^{-7}$$

gleich Wahrscheinlich $\Rightarrow p = 0,5$

$$\frac{1}{0,5} = 1 + e^{(5+x)} \quad | -1$$

$$1 = e^{(5+x)} \quad |\ln$$

$$\ln 1 = \ln e^{(5+x)}$$

$$0 = 5 + x \quad | -5$$

$$-5 = x$$

3. Funktion knn.cv

sie führt eine Kreuzvalidierung an dem
gegebenen Datensatz durch

mit euklidischer Distanz