

# Hausaufgabe 9

## ① Trainingsdaten

	A	B	C		D	E	F	G	H
x	2	1	2	x	6	6	5	4	3
y	8	7	6	y	4	3	2	2	2
z	5	5	5	z	5	5	5	5	5

krank

gesund

Beachte, dass  $z=5$  für alle Werte. Damit vereinfacht sich das Problem auf ein zweidimensionales:

$$\alpha = (1, 6)^T: \text{Distanzen } \left. \begin{array}{l} d_{\alpha A} = \sqrt{5} \\ d_{\alpha B} = 1 \\ d_{\alpha C} = 1 \end{array} \right\} \begin{array}{l} d_{\alpha D} = \sqrt{29} \\ d_{\alpha E} = \sqrt{34} \\ d_{\alpha F} = \sqrt{32} \\ d_{\alpha G} = 5 \\ d_{\alpha H} = \sqrt{20} \end{array} \left. \begin{array}{l} 1. kNN(\alpha, k=1) = \{B\} = \{C\} \\ 2. kNN(\alpha, k=3) = \{A, B, C\} \\ 3. kNN(\alpha, k=5) = \{A, B, C, H, G\} \end{array} \right\} \begin{array}{l} \text{krank} \\ \text{krank} \\ \text{krank} \end{array}$$

$\alpha$  stammt mit hoher Wahrscheinlichkeit aus krankem Gewebe

Für  $k=7$  ist die Entscheidung immer gesund, da um 3 Datensätze krank gelabelt sind. Deswegen hat die Berechnung dieser Menge keinen Sinn.

$$\beta = (5, 3)^T: \left. \begin{array}{l} d_{\beta A} = \sqrt{34} \\ d_{\beta B} = \sqrt{32} \\ d_{\beta C} = \sqrt{18} \end{array} \right\} \left. \begin{array}{l} d_{\beta D} = \sqrt{2} \\ d_{\beta E} = 1 \\ d_{\beta F} = 1 \\ d_{\beta G} = \sqrt{2} \\ d_{\beta H} = \sqrt{5} \end{array} \right\} \begin{array}{l} 1. kNN(\beta, k=1) = \{E\} = \{F\} \\ 2. kNN(\beta, k=3) = \{E, F, D(G)\} \\ 3. kNN(\beta, k=5) = \{D, E, F, G, H\} \end{array} \begin{array}{l} \text{gesund} \\ \text{gesund} \\ \text{gesund} \end{array}$$

$\beta$  stammt mit hoher Wahrscheinlichkeit aus gesundem Gewebe

$$\gamma = (2, 4)^T: \left. \begin{array}{l} d_{\gamma A} = 4 \\ d_{\gamma B} = \sqrt{10} \\ d_{\gamma C} = 2 \end{array} \right\} \left. \begin{array}{l} d_{\gamma D} = 4 \\ d_{\gamma E} = \sqrt{17} \\ d_{\gamma F} = \sqrt{13} \\ d_{\gamma G} = \sqrt{8} \\ d_{\gamma H} = \sqrt{5} \end{array} \right\} \begin{array}{l} 1. kNN(\gamma, k=1) = \{C\} \\ 2. kNN(\gamma, k=3) = \{C, H, G\} \\ 3. kNN(\gamma, k=5) = \{C, H, G, B, F\} \end{array} \begin{array}{l} \text{krank} \\ \text{gesund} \\ \text{gesund} \end{array}$$

Ein Zwischung von  $\gamma$  kann nur unter Vorbehalt vorgenommen werden, da die kNN Methode für unterschiedliches  $k$  andere Ergebnisse produziert.

## Aufgabe 9.2 Logistische Regression (30%)

Sie betrachten die Wahrscheinlichkeiten  $p(x)$ , dass ein Patient erkrankt ist. Hierfür haben Sie für verschiedene Patienten die Expression  $x$  eines bestimmten Genes ermittelt. Eine logistische Regression der funktionalen Form

$$p(x) = \frac{1}{1 + \exp(-\beta_0 + \beta_1 x)} \quad (4)$$

hat die Parameter  $\beta_0 = -5$  und  $\beta_1 = 1$  ergeben.

decision boundary bei  $p(x_{db}) = 0,5$ . Umstellen nach  $x$ :

$$\frac{1}{p} = 1 + \exp(-\beta_0 + \beta_1 x)$$

$$\Leftrightarrow \exp(-\beta_0 + \beta_1 x) = \frac{1}{p} - 1$$

$$\Leftrightarrow \beta_1 x = \ln\left(\frac{1}{p} - 1\right) + \beta_0$$

$$x = \frac{\ln\left(\frac{1}{p} - 1\right) + \beta_0}{\beta_1} \quad \text{und mit } p=0,5; \beta_0=-5 \text{ und } \beta_1=1 \text{ ist}$$

$$\underline{\underline{x_{db} = -5}}$$

$$\text{Weiterhin ist } p(x=2) = 0,911 \cdot 10^{-3} \approx 0,9 \%$$

$$p(x=6) = 0,017 \cdot 10^{-3} \approx 0,02 \%$$

$$p(x=10) = 0,0003 \cdot 10^{-3} \approx 0,0003 \%$$