

# Ethik in der KI

KI soll **verantwortlich** sein - idealerweise sollten wir wissen, wer oder was von einem KI-Modell getroffene Entscheidung verursacht hat.

KI soll **Fair** sein - Frei von Vorurteilen.

KI soll **Mensch-zentriert** sein - Sie sollte menschliche Bedürfnisse in sozialen, politischen und emotionalen Kontexten beachten und respektieren.

→ Wir können diese Anforderungen zu erfüllen, indem wir KI erklärbar machen.

→ Ein möglicher Ansatz ist die Verwendung von symbolischer KI.

- Ein besonderer Ansatz ist die Verwendung von symbolischen Vektorarchitekturen, auch bekannt als 'hyperdimensional computing' (HDC)

# HDC

Binding:

$$\text{Wein} = \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix} \quad \text{Rot} = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} \quad \text{Weiß} = \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix}$$

↙ Element-wise multiplication

$$\text{Rotwein} = \text{Wein} \otimes \text{Rot}$$

$$= \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix} \otimes \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix}$$

$$\text{Weißwein} = \text{Wein} \otimes \text{Weiß}$$

$$= \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix} \otimes \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix}$$

Bundling:

Gegeben - eine Menge von HDVs  $\{\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n\}$

Use majority vote sum to bundle them:

$$S = \sum_{k=1}^n v_k$$

Then take the sign of  $S$ :

$$\text{Sign}(s) = \begin{cases} +1 & s \geq 0 \\ -1 & s < 0 \end{cases}$$

Weinkeller = Rotwein  $\oplus$  Weißwein

$$\Rightarrow s = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} \oplus \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ -2 \\ -2 \end{bmatrix}$$

$$\text{Weinkeller} = \text{sign}(s) = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix}$$

## Unbinding

Problem: Wir haben ein unbekanntes Getränk mit einer rosa Farbe - eine Mischung aus Rot und Weiß.

Lösung: Wir können den gebündelten Vektor Weinkeller abfragen ob es Wein ist.

$$\text{Rosa} = \text{Rot} \oplus \text{Weiß}$$

$$= \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} \oplus \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ 2 \\ -2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}$$

Unbinding:

$$\text{Rosa} \otimes \text{Weinkeller} = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} \otimes \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}$$

Ergebniss  $\neq$  Wein

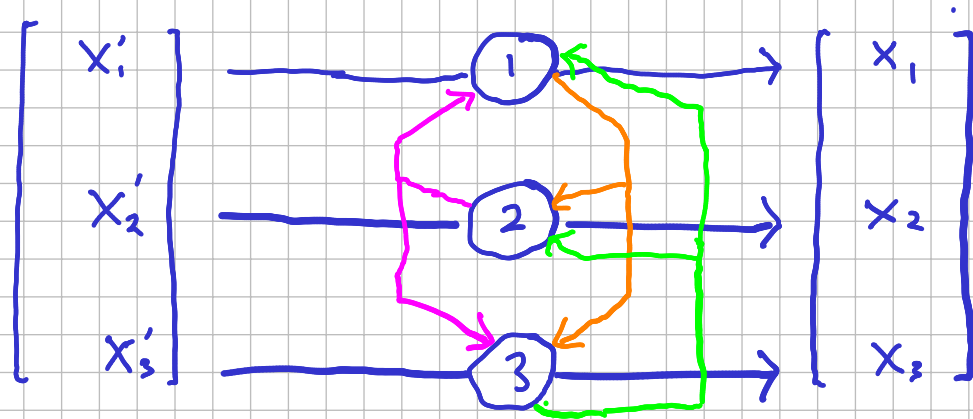
Ergebniss ist eine verrauschte Version von Wein

Frage: Wie können wir die ursprüngliche Repräsentation von Wein wiederherstellen?

# Hopfield Network (HF)

Ein Hopfield-Netz ist ein vollständig verbundenes neuronales Netz, das einen verrauschten oder unvollständigen Vektor aufnimmt und die am besten passende gelernte Repräsentation ausgibt.

$\vec{x}' \Rightarrow$  verrauschten  $\vec{x}$



- No 'self-loops'

- Gewichte  $W \in \mathbb{R}^{N \times N} \Rightarrow$  hier  $3 \times 3$

$$W = \begin{bmatrix} 0 & w_{12} & w_{13} \\ w_{21} & 0 & w_{23} \\ w_{31} & w_{32} & 0 \end{bmatrix}$$

Color-coded arrows point from the weights to the connections in the diagram above: orange for  $w_{12}$  and  $w_{21}$ , pink for  $w_{13}$  and  $w_{31}$ , and green for  $w_{23}$  and  $w_{32}$ .

$$w_{ij} = \frac{1}{N} \sum_{\mu=1}^P x_i^{(\mu)} x_j^{(\mu)}$$

Beispiel: Wir können die Muster von Wein, Rot und Weiß lernen.

$$P = 3$$

$$N = 3$$

$$p_1 = \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix}$$

Wein

$$p_2 = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}$$

Rot

$$p_3 = \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix}$$

Weiß

$$w_{12} = \frac{1}{3} [ (-1)(-1) + (1)(1) + (-1)(1) ] = \frac{1}{3} (1) = \frac{1}{3}$$

$$w_{13} = \frac{1}{3} [ (-1)(1) + (1)(-1) + (-1)(-1) ] = \frac{1}{3} (-1) = -\frac{1}{3}$$

$$w_{21} = w_{12} = \frac{1}{3}$$

$$w_{23} = \frac{1}{3} [ (-1)(1) + (1)(-1) + (1)(-1) ] = \frac{1}{3} (-3) = -1$$

$$w_{31} = w_{13} = -\frac{1}{3}$$

$$w_{32} = w_{23} = -1$$

$$W = \begin{bmatrix} 0 & \frac{1}{3} & -\frac{1}{3} \\ \frac{1}{3} & 0 & -1 \\ -\frac{1}{3} & -1 & 0 \end{bmatrix}$$

Frage: Wie können wir die ursprüngliche Repräsentation von Wein wiederherstellen?

$\vec{X}$  ausgabe für  $\vec{X}'$  eingabe

$$X_i = \text{sign} \left( \sum_j w_{ij} X'_j \right)$$

$$\text{Wein} + \text{Noise (Rauschen)} = \vec{X}' = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}$$

$$\begin{aligned} X_1 &= \text{sign} [w_{12} X'_2 + w_{13} X'_3] = \text{sign} \left[ \frac{1}{3}(-1) + \left(-\frac{1}{3}\right)(1) \right] \\ &= \text{sign} \left( -\frac{2}{3} \right) = \textcircled{-1} \end{aligned}$$

$$\begin{aligned} X_2 &= \text{sign} [w_{21} X'_1 + w_{23} X'_3] = \text{sign} \left[ \frac{1}{3}(1) + (-1)(1) \right] \\ &= \text{sign} \left( -\frac{2}{3} \right) = \textcircled{-1} \end{aligned}$$

$$\begin{aligned} X_3 &= \text{sign} [w_{31} X'_1 + w_{32} X'_2] = \text{sign} \left[ \left(-\frac{1}{3}\right)(1) + (-1)(-1) \right] \\ &= \text{sign} \left( \frac{2}{3} \right) = \textcircled{1} \end{aligned}$$

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix} = \text{Wein} \quad \checkmark$$