

Support Vector Machines & K-Nearest Neighbor



Universität Stuttgart



Daniel Fink

Quantum Machine Learning – Seminar

24. Januar 2020



Gesichtserkennung

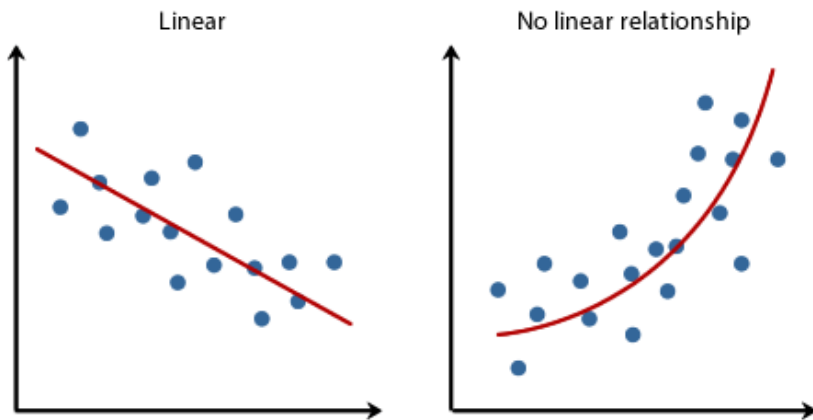
{ → Face Detection
→ Face Recognition

Here's an example of my normal handwriting. Notice how letters are connected in a way that has little to do with word length - as in "normal" - 1, 2, 3 letters are grouped.

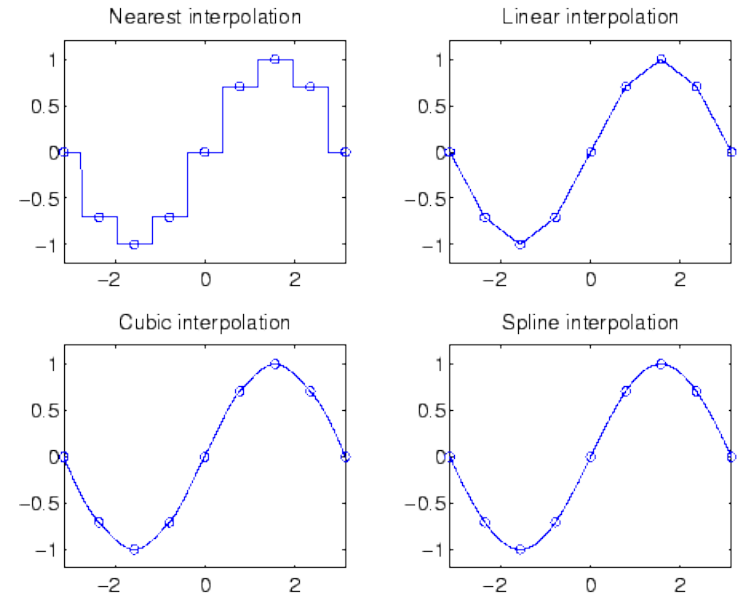


Here's an example of my normal handwriting. Notice how letters are connected in a way that has little to do with word length - as in "normal" - 1, 2, 3 letters are grouped.

Handschrifterkennung



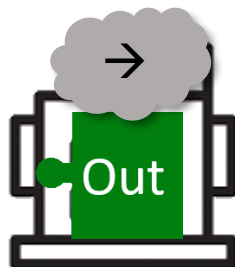
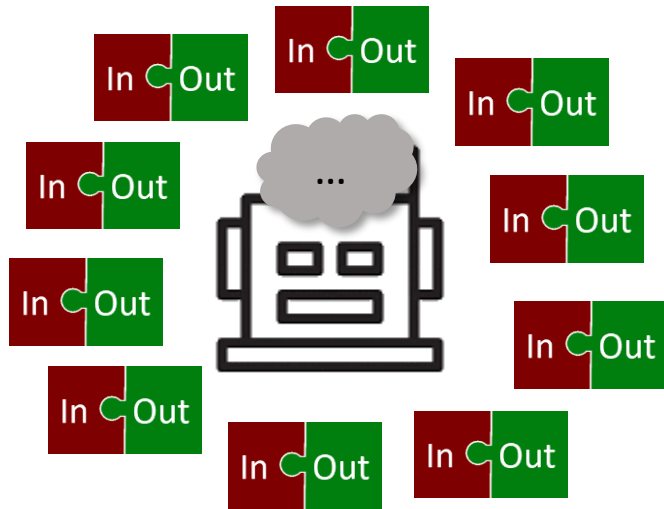
Regression



Interpolation

Supervised Learning





K-Nearest Neighbor

Theorie

Auf Quanten-Computer

Support Vector Machines

Theorie

Im Kontext von NISQ

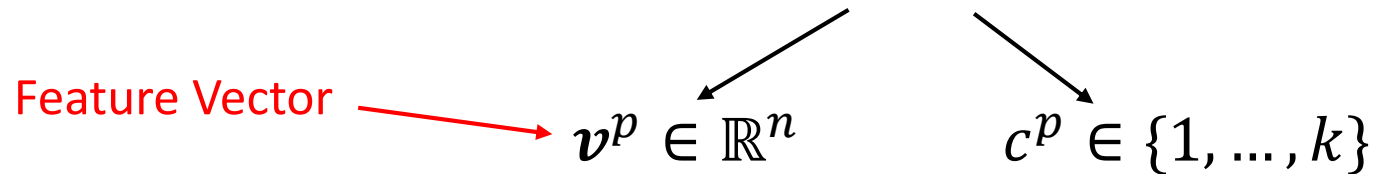
Für Big Data



K-Nearest Neighbor - Theorie

K-Klassifizierungsproblem:

- Gegeben: Trainingsmenge $T = \{(\mathbf{v}^p, c^p)\}_{p=1, \dots, N}$

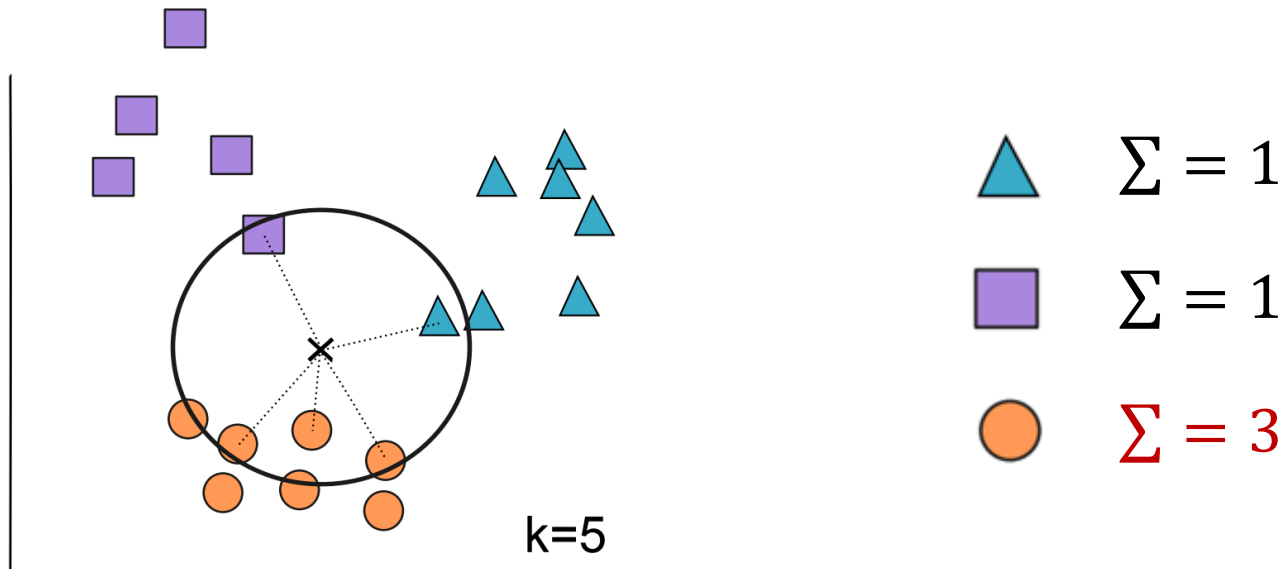


Zu klassifizierender Input-Vektor $\tilde{\mathbf{x}} \in \mathbb{R}^n$

- Gesucht: Klassifizierung $\tilde{c} \in \{1, \dots, k\}$

Idee:

- Wähle die k nächsten Nachbarn von \tilde{x}
- Weise \tilde{x} die Klasse der Mehrheit der Nachbarn zu



Variationen:

- *K-Nearest Neighbor* mit *Distanz-Gewichtung*
 - Gewichte die Nachbarn mit $1/dist$
- *K-Nearest Centroids*
 - Berechne vorab die jeweiligen Zentren der Klassen
 - Wähle die Klasse mit dem nächsten Zentrum

Frage:

- Was bedeutet „am nächsten“?
- Geeignetes Ähnlichkeits- bzw. Distanzmaß notwendig

Distanzmaß	$d(i, j)$
Manhattan-Metrik (L_1)	$\sum_k x_{i,k} - x_{j,k} $
Euklidische Metrik (L_2)	$\sqrt{\sum_k x_{i,k} - x_{j,k} ^2}$
Hamming-Metrik	$ \{k \in \{1, \dots, n\} : x_{i,k} \neq x_{j,k}\} $
...	...



K-Nearest Neighbor – Auf Quanten-Computer

K-Quanten-Klassifizierungsproblem:

- Gegeben: Trainingsmenge

$$\{ |v_1^p \dots v_n^p, c^p\rangle \}_{p=1, \dots, N} \subset H_2^{\otimes n} \otimes H_k$$

k Klassen
↙

Zu klassifizierender Input-Vektor $|\tilde{x}_1 \dots \tilde{x}_n\rangle \in H_2^{\otimes n}$

- Gesucht: Klassifizierung $|\tilde{c}\rangle \in H_k$

→ Kodierung der Daten als Bitstrings

→ Verwende Hamming-Distanz

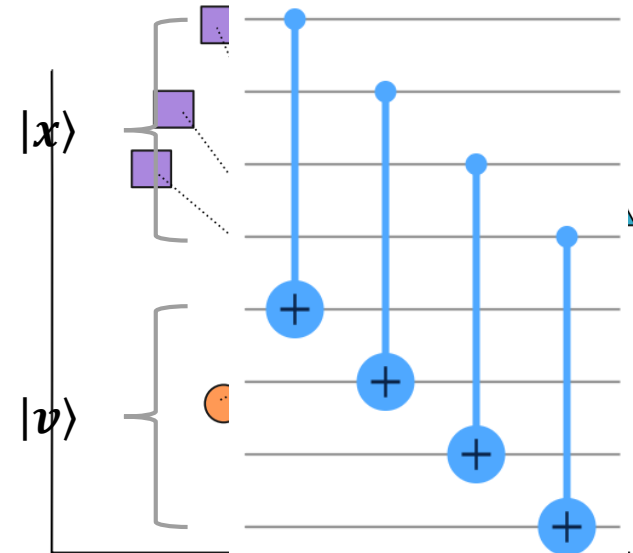
Idee:

- K-Nearest Neighbor mit Distanzgewichtung
- Betrachte alle Trainingsdaten
- Berechne Hamming-Distanz mit CNOT

$$x = 0\mathbf{1}101\mathbf{1}01$$

$$v = 00101011$$

➔ $d(x, v) = 3$



$$|x\rangle = |0\mathbf{1}101\mathbf{1}01\rangle$$

$$|v\rangle = |00101011\rangle$$

CNOT



$$|x'\rangle = |01101101\rangle$$

$$|v'\rangle = |0\mathbf{1}00\mathbf{0}110\rangle \rightarrow \Sigma = 3$$

Ablauf:

1) Erzeuge Superposition der Trainingsdaten

$$|T\rangle = \sum_p |v_1^p \dots v_n^p, c^p\rangle$$

2) Superpositioniere mit Input-Zustand

$$|\psi_o\rangle = \sum_p |\tilde{x}_1 \dots \tilde{x}_n; v_1^p \dots v_n^p, c^p; 0\rangle$$

Input

Training

Anchilla

3) Führe Hadamard auf Anchilla aus

$$|\psi_1\rangle = \sum_p |\tilde{x}_1 \dots \tilde{x}_n; v_1^p \dots v_n^p, c^p\rangle \otimes (|0\rangle + |1\rangle)$$

Ablauf:

4) Berechne Hamming-Distanz mit CNOT

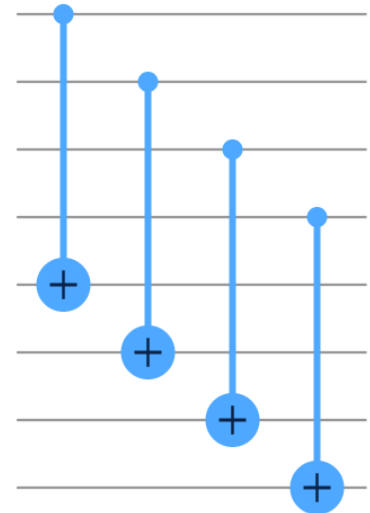
$$|\psi_2\rangle = \prod_k CNOT(\tilde{x}_k, v_k^p) |\psi_1\rangle$$

- Derzeit:

Große Ähnlichkeit = **Kleine** Hamming-Distanz

- Für später sinnvoller:

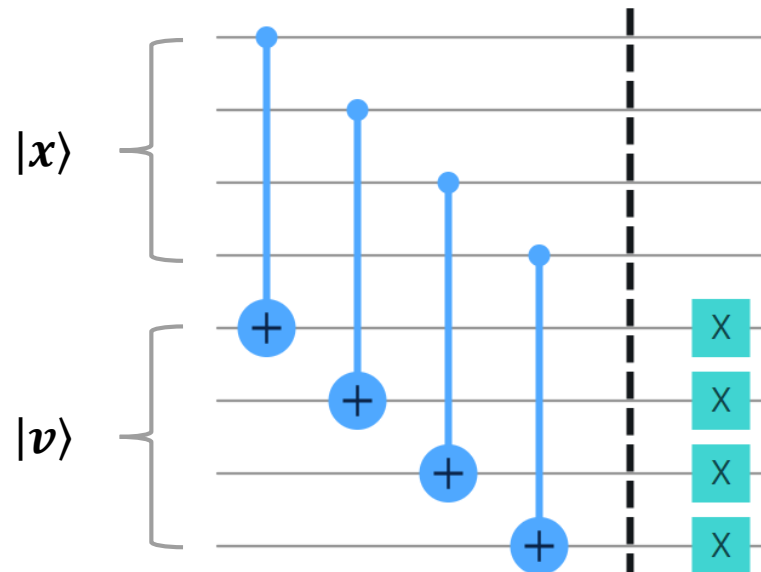
Große Ähnlichkeit = **Große inverse** Hamming-Distanz



Ablauf:

4) Berechne **inverse** Hamming-Distanz mit CNOT

$$|\psi_2\rangle = \prod_k X(\tilde{x}_k) CNOT(\tilde{x}_k, v_k^p) |\psi_1\rangle$$



Ablauf:

4) Berechne **inverse** Hamming-Distanz mit CNOT

$$\begin{aligned} |\psi_2\rangle &= \prod_k X(\tilde{\mathbf{x}}_k) CNOT(\tilde{\mathbf{x}}_k, \mathbf{v}_k^p) |\psi_1\rangle \\ &= \sum_p |\tilde{\mathbf{x}}_1 \dots \tilde{\mathbf{x}}_n; d_1^p \dots d_n^p, c^p\rangle \otimes (|0\rangle + |1\rangle) \end{aligned}$$

5) Summiere **inverse Hamming-Distanzen**

$$U = e^{i\frac{\pi}{2n}H} \text{ mit } H = 1 \otimes \underbrace{\sum_k \left(\frac{\sigma_z+1}{2}\right)_k}_{\text{Summation}} \otimes 1 \otimes \sigma_z$$

Negativ, wenn
Ancilla = 1

$$\left(\frac{\sigma_z+1}{2}\right) |0\rangle = \frac{1 \cdot |0\rangle + |0\rangle}{2} = |0\rangle$$

$$\left(\frac{\sigma_z+1}{2}\right) |1\rangle = \frac{-1 \cdot |1\rangle + |1\rangle}{2} = 0$$

Ablauf:

5) Summiere inverse Hamming-Distanzen

$$\begin{aligned} |\psi_3\rangle &= U |\psi_2\rangle \\ &= \sum_p e^{i\frac{\pi}{2n}d_H(\tilde{x}, v^p)} |\tilde{x}_1 \dots \tilde{x}_n; d_1^p \dots d_n^p, c^p; \mathbf{0}\rangle + e^{-i\frac{\pi}{2n}d_H(\tilde{x}, v^p)} |\tilde{x}_1 \dots \tilde{x}_n; d_1^p \dots d_n^p, c^p; \mathbf{1}\rangle \end{aligned}$$

6) Wende erneut Hadamard auf Anchilla an

$$\begin{aligned} |\psi_4\rangle &= (1 \otimes 1 \otimes 1 \otimes H) |\psi_3\rangle \\ &= \sum_p \cos\left[\frac{\pi}{2n}d_h(\tilde{x}, v^p)\right] |\tilde{x}; d^p, c^p; 0\rangle + i \sin\left[\frac{\pi}{2n}d_h(\tilde{x}, v^p)\right] |\tilde{x}; d^p, c^p; 1\rangle \end{aligned}$$

Mit Eulerscher Formel: $e^{ix} = \cos(x) + i \sin(x)$

Ablauf:

$$\cos(0) = 1 \quad \sin(0) = 0$$

7) Messe Anchilla

$$P(|a\rangle = |0\rangle) \propto \sum_p \cos^2 \left[\frac{\pi}{2n} d_h(\tilde{x}, v^p) \right] \Rightarrow \text{Input weit weg von Trainingsdaten}$$

$$P(|a\rangle = |1\rangle) \propto \sum_p \sin^2 \left[\frac{\pi}{2n} d_h(\tilde{x}, v^p) \right] \Rightarrow \text{Input nah dran an Trainingsdaten}$$

$$d_h(\tilde{x}, v^p) \begin{cases} \text{Klein, wenn } \tilde{x} \text{ und } v^p \text{ unähnlich sind} \\ \text{Groß, wenn } \tilde{x} \text{ und } v^p \text{ ähnlich sind} \end{cases}$$

8) Verwende Threshold bis Anchilla = 1 gemessen wird

$$\rightarrow \text{Messe } m \in \mathbb{N} \text{ mal} \begin{cases} \forall m: \text{Anchilla} = 0 \rightarrow \text{Abbruch!} \\ \exists m: \text{Anchilla} = 1 \rightarrow \text{Bestimme Klasse } \tilde{c} \end{cases}$$

Ablauf:

8) Verwende Threshold bis **Anchilla = 1** gemessen wird

$$\begin{aligned} |\psi_5\rangle &= \sum_p i \sin \left[\frac{\pi}{2n} d_h(\tilde{\mathbf{x}}, \mathbf{v}^p) \right] |\tilde{\mathbf{x}}; \mathbf{d}^p, c^p; 1\rangle \\ &= \sum_c |c\rangle \otimes \sum_{l \in c} i \sin \left[\frac{\pi}{2n} d_h(\tilde{\mathbf{x}}, \mathbf{v}^l) \right] |\tilde{\mathbf{x}}; \mathbf{d}^l; 1\rangle \end{aligned}$$

9) Messe Klassen-Qbit

$$P(c) \propto \sum_{l \in c} \sin^2 \left[\frac{\pi}{2n} d_h(\tilde{\mathbf{x}}, \mathbf{v}^l) \right]$$

Summe läuft über die
Trainingsdaten der Klasse c

Proportional zur inversen
Hamming-Distanz

Fazit:

- Konstruktion einer beliebigen Superposition $\sum_p |\tilde{x}_1 \dots \tilde{x}_n; v_1^p \dots v_n^p, c^p; 0\rangle$ ist schwierig \rightarrow QRAM
- Laufzeit $O(TNn)^4$
- Qbits $O(n)$
- Liefert mehr Informationen \rightarrow WS für die Klassen





Support Vector Machine - Theorie

Binäres Klassifizierungsproblem:

- Gegeben: Trainingsmenge $T = \{(\mathbf{v}^p, c^p)\}_{p=1, \dots, N}$

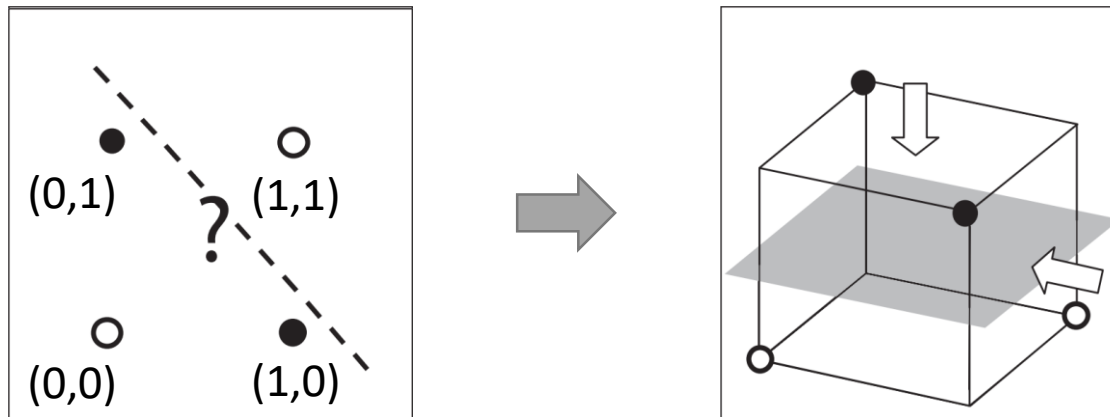
$$\begin{array}{ccc} & \swarrow & \searrow \\ \mathbf{v}^p \in \mathbb{R}^n & & c^p \in \{-1, +1\} \end{array}$$

Zu klassifizierender Input-Vektor $\tilde{\mathbf{x}} \in \mathbb{R}^n$

- Gesucht: Klassifizierung $\tilde{c} \in \{-1, +1\}$

Idee:

- Transformiere Daten in einen höherdimensionalen Raum
- Finde Hyperebene, welche die Datenpunkte trennt
- Prüfe auf welcher Seite \tilde{x} liegt und weise Klasse zu
- Beispiel XOR-Problem:



Geometrische Herleitung:

- Hyperebene als Punkt-Normalenform

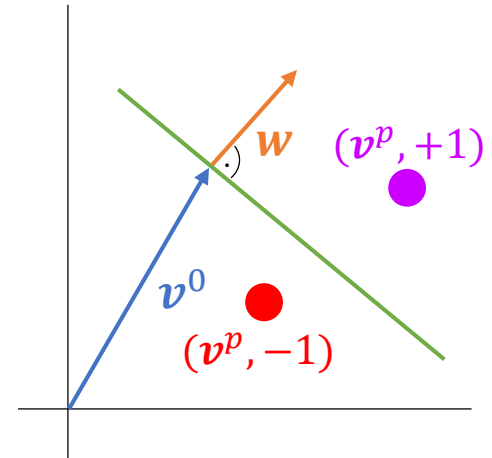
$$\mathbf{w}^T(\mathbf{v} - \mathbf{v}^0) = 0$$

$$\mathbf{w}^T \mathbf{v} + b = 0 \text{ mit } b = -\mathbf{w}^T \mathbf{v}^0$$

- Trennung der Trainingsdaten durch

$$\mathbf{w}^T \mathbf{v}^p + b \geq 0 \quad \text{für alle } c^p = +1$$

$$\mathbf{w}^T \mathbf{v}^p + b < 0 \quad \text{für alle } c^p = -1$$

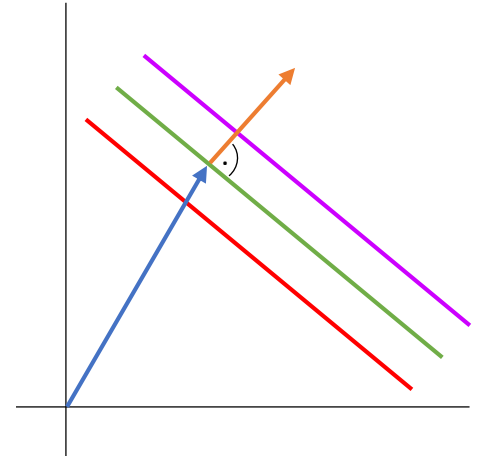


Geometrische Herleitung:

- Führe Trennbereich ein

$$\mathbf{w}^T \mathbf{v}^p + \mathbf{b} \geq +1 \quad \text{für alle } c^p = +1$$

$$\mathbf{w}^T \mathbf{v}^p + \mathbf{b} < -1 \quad \text{für alle } c^p = -1$$



- Damit gilt

$$c^p (\mathbf{w}^T \mathbf{v}^p + \mathbf{b}) \geq +1$$

variabel

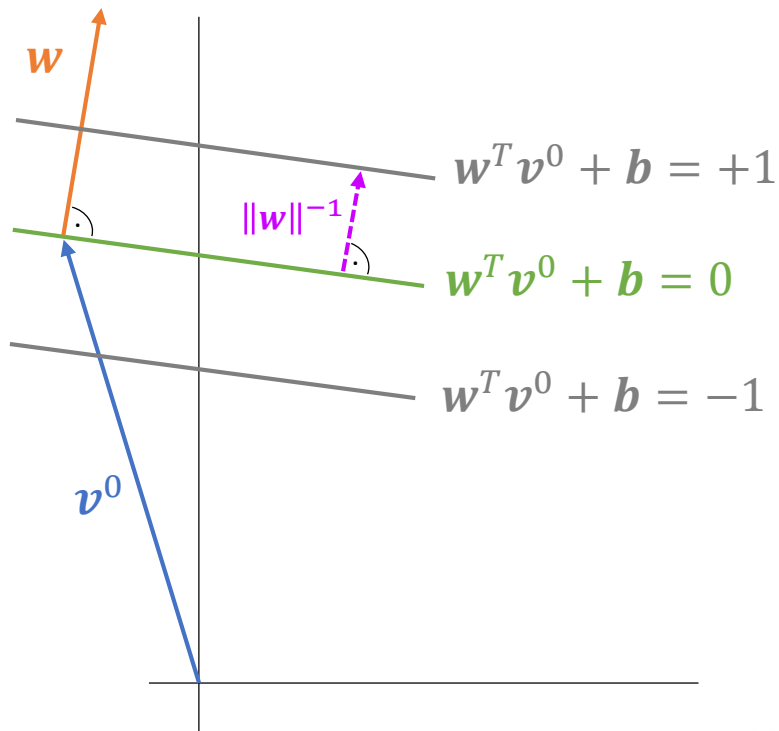
b hat keinen Einfluss auf die Breite



Bestimme \mathbf{w} so, dass diese Gleichung erfüllt wird

Geometrische Herleitung:

$$c^p(\mathbf{w}^T \mathbf{v}^p - b) \geq +1$$



Trennbereich soll möglichst groß werden

→ Minimiere Kostenfunktion

$$\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$$

Bisher:

- Minimiere

$$\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$$

Maximaler Trennbereich

- Erfülle

$$c^p (\mathbf{w}^T \mathbf{v}^p - b) \geq +1$$

Trennbarkeit

→ Quadratisches Optimierungsproblem

→ Verwende Lagrange-Multiplikatoren

Minimiere

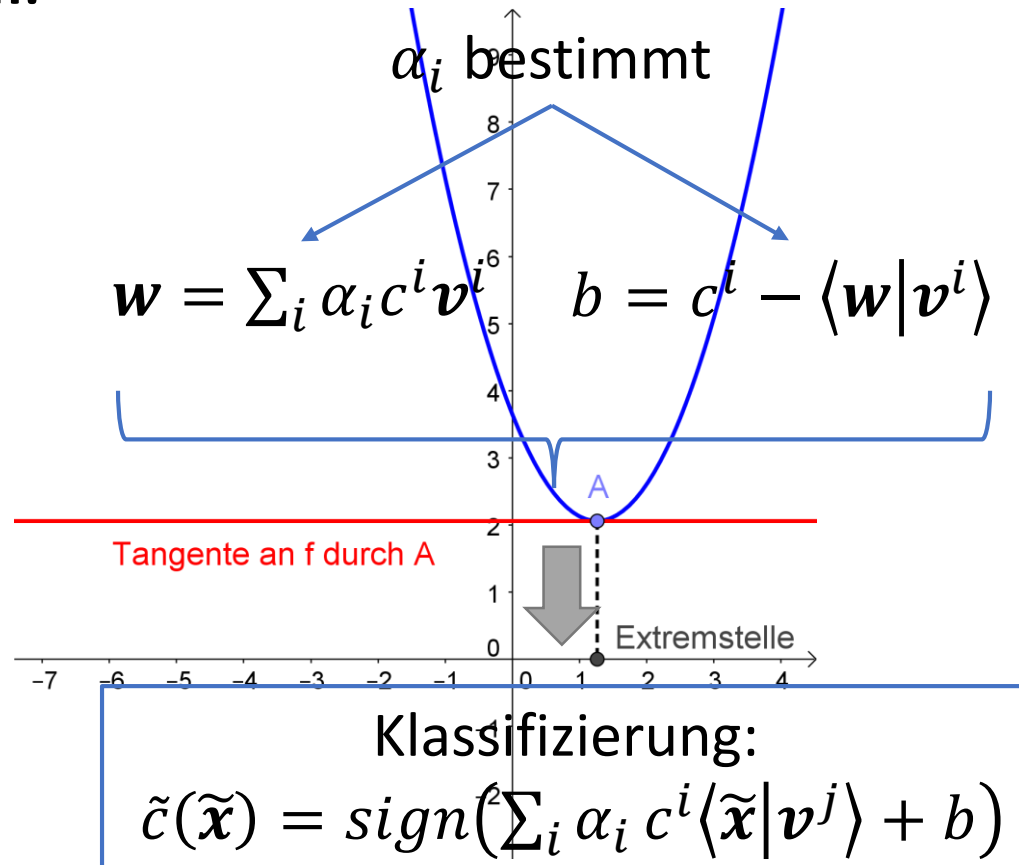
$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_p \alpha_p (c^p (\mathbf{w}^T \mathbf{v}^p - b) - 1)$$

hinsichtlich \mathbf{w} und b

Optimierungsproblem:

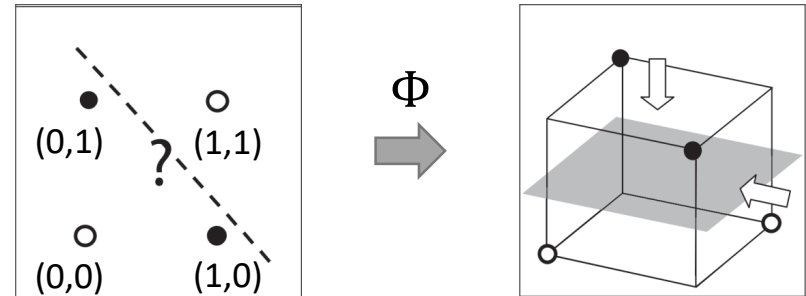
$$\frac{\partial \mathcal{L}(\mathbf{w}, b, \alpha)}{\partial \mathbf{w}} = 0$$

$$\frac{\partial \mathcal{L}(\mathbf{w}, b, \alpha)}{\partial b} = 0$$



Bisher:

- Daten wurden noch nicht in höherdimensionalen Raum transformiert
- Verwende Kernel-Trick



$$\Phi: \mathbb{R}^n \rightarrow H$$

nicht gut trennbar

mehr Freiheiten für Trennung

$$\langle x|y \rangle \rightarrow \underbrace{\langle \Phi(x)|\Phi(y) \rangle}_{\text{schwer zu berechnen}} = \underbrace{K(x,y)}_{\text{einfach zu berechnen}}$$

schwer zu berechnen

einfach zu berechnen



Support Vector Machine - NISQ

Vorher:

- Klassischer Algorithmus → Quanten Algorithmus

Jetzt:

- Quanten-Rechner → Quanten Klassifikator

Konkret:

- IBMQ mit 5 Qbits
- Keine Fehlerkorrektur
- Eingeschränkte Konnektivität
- Eingeschränkte Breite
- Eingeschränkte # Operatoren



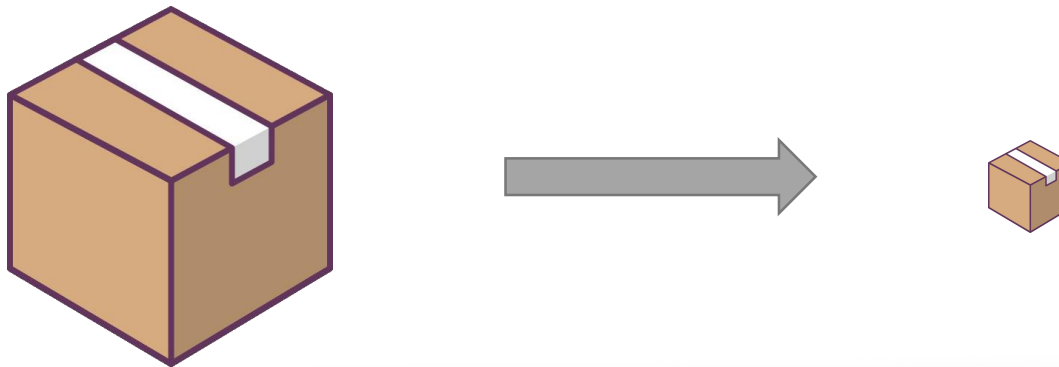
Was kann man
damit anstellen?

Kodierung der Daten:

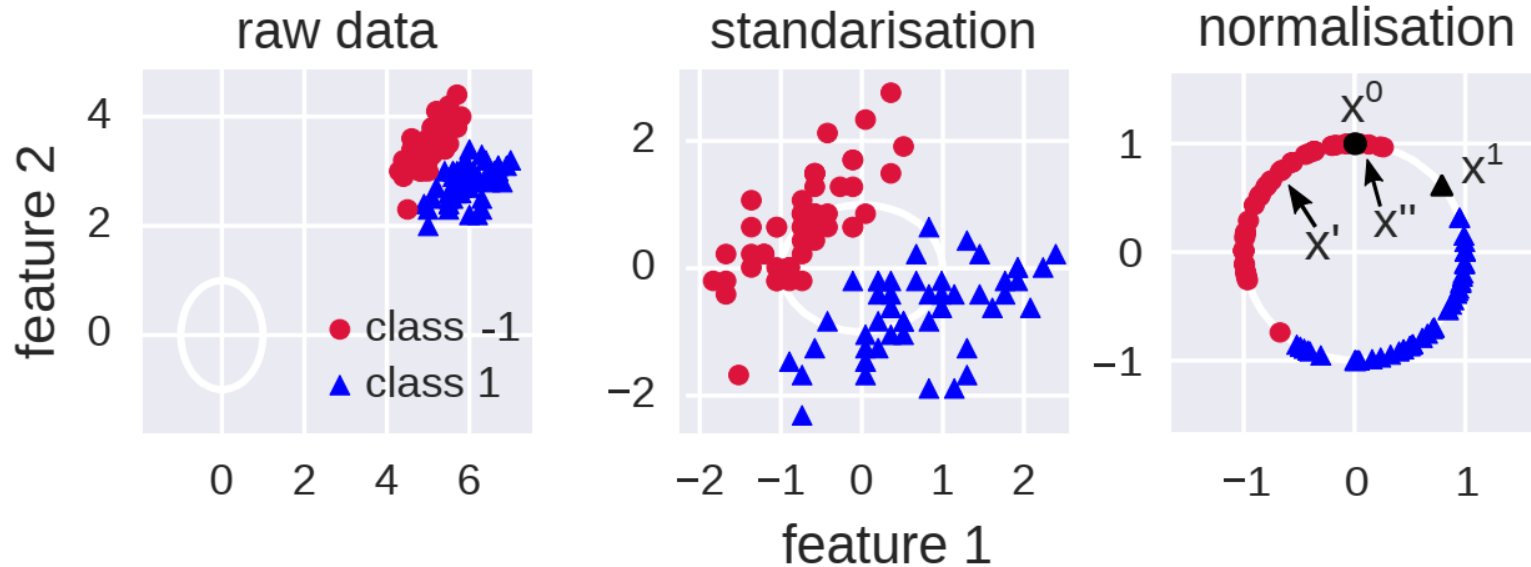
$$\bullet \quad \underbrace{v = (v_1, v_2, \dots, v_n)}_{O(n) \text{ Bits}} \in \mathbb{R}^n \rightarrow |\Psi_v\rangle = \underbrace{\sum_{i=0}^{n-1} v_i |i\rangle}_{n \text{ Terme}} \in H_2^{\otimes d}$$

→ Superposition: $n = 2^d \rightarrow d = \log_2(n)$

→ Exponentielle „Komprimierung“ der Daten



(Klassische) Aufbereitung der Daten:



Ablauf:

1) Superpositioniere und verschränke Input und Trainingsdaten

$$|\Psi_0\rangle = \sum_p |p\rangle (|0\rangle |\Psi_{\tilde{x}}\rangle + |1\rangle |\Psi_{v^p}\rangle) |c^p\rangle$$

Kennzeichnung Anchilla Input/Training Klasse

2) Führe Hadamard auf Anchilla aus

$$|\Psi_1\rangle = \sum_p |p\rangle (|0\rangle |\Psi_{\tilde{x}+v^p}\rangle + |1\rangle |\Psi_{\tilde{x}-v^p}\rangle) |c^p\rangle$$

Kennzeichnung Anchilla Input/Training Klasse

$$|\Psi_{\tilde{x} \pm v^p}\rangle = |\Psi_{\tilde{x}}\rangle \pm |\Psi_{v^p}\rangle$$

Ablauf:

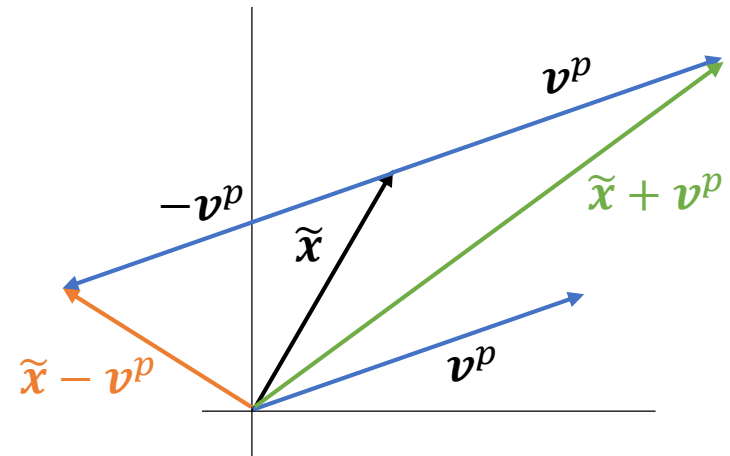
$$|\Psi_1\rangle = \sum_p |p\rangle (|0\rangle |\Psi_{\tilde{x}+v^p}\rangle + |1\rangle |\Psi_{\tilde{x}-v^p}\rangle) |c^p\rangle$$

3) Messe Anchilla

$$\begin{aligned} |\Psi_{\tilde{x}}\rangle + |\Psi_{v^p}\rangle &= \sum_{i=0}^{n-1} \tilde{x}_i |i\rangle + \sum_{i=0}^{n-1} v_i^p |i\rangle \\ &= \sum_{i=0}^{n-1} (\tilde{x}_i + v_i^p) |i\rangle \end{aligned}$$

$$\rightarrow P(|a\rangle = |0\rangle) \propto \sum_p |\tilde{x} + v^p|^2$$

→ Hohe WS wenn Input nah an Trainingsdaten



Ablauf:

3) Resultierender Zustand

$$|\Psi_2\rangle = \sum_p \sum_{i=0}^{n-1} |p\rangle \underbrace{(\tilde{x}_i + v_i^p)}_{\text{Gewichtung von Klasse } c^p} |i\rangle |c^p\rangle$$

Gewichtung von Klasse c^p

4) Messe Klasse

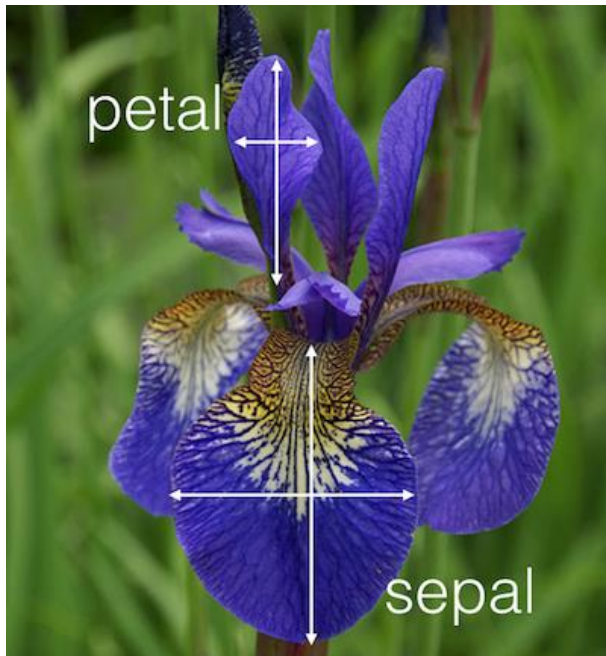
$$P(|c\rangle = |0\rangle) \propto \sum_{p:c^p=0} |\tilde{x} + v^p|^2$$

$$\Rightarrow \tilde{c} = -1$$

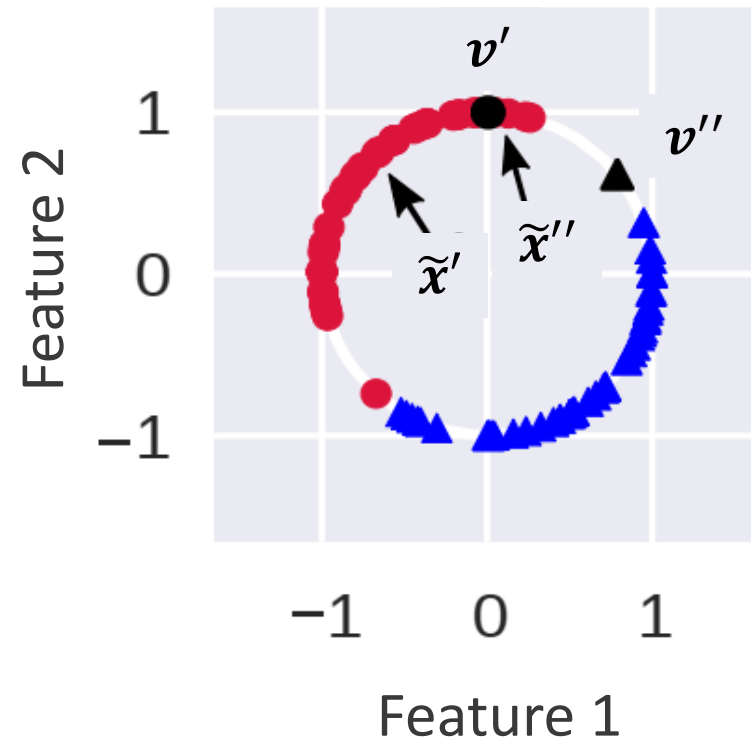
Klassifizierung: $\tilde{c}(\tilde{x}) = \text{sign}\left(\sum_p c^p \left[1 - \frac{1}{4N} |\tilde{x} - v^p|^2\right]\right)$

Kernel

Live Demo:

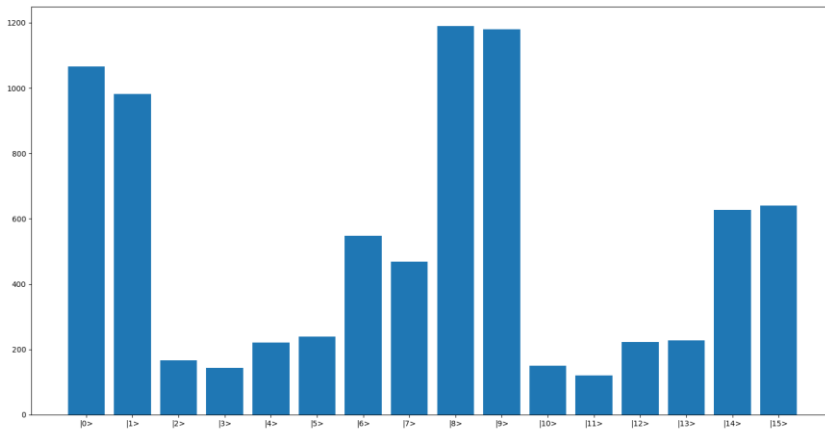


● = -1 ▲ = +1



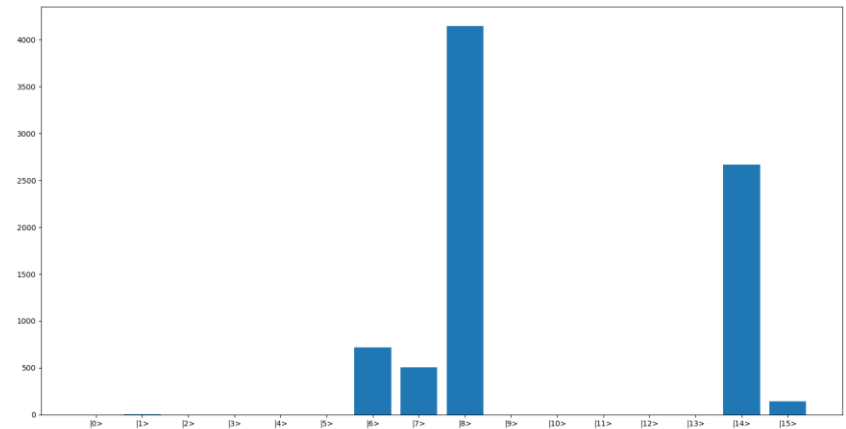


Live Demo:



IBMQ

$$P(\tilde{c} = -1) \approx 60\%$$



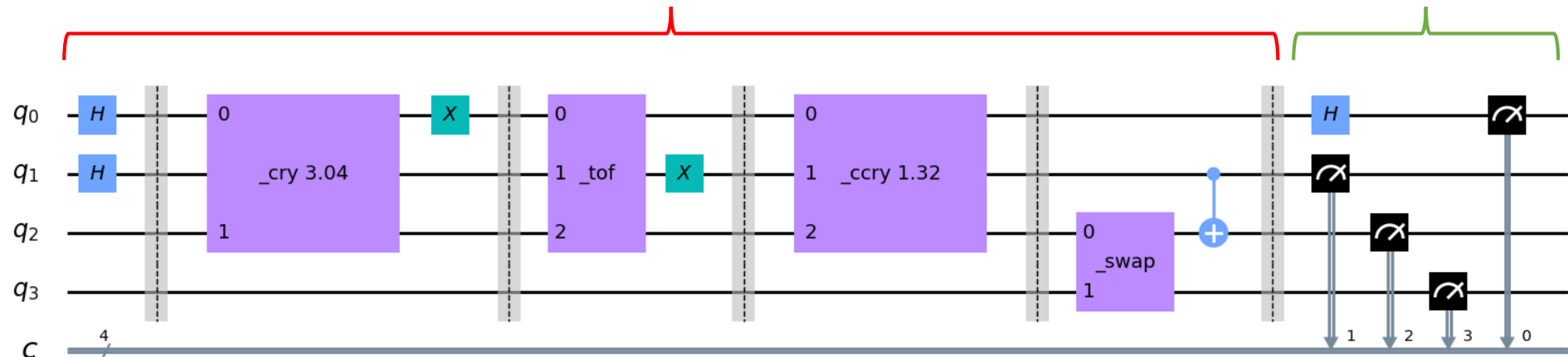
Simulator

$$P(\tilde{c} = -1) = 55\%$$

Live Demo:

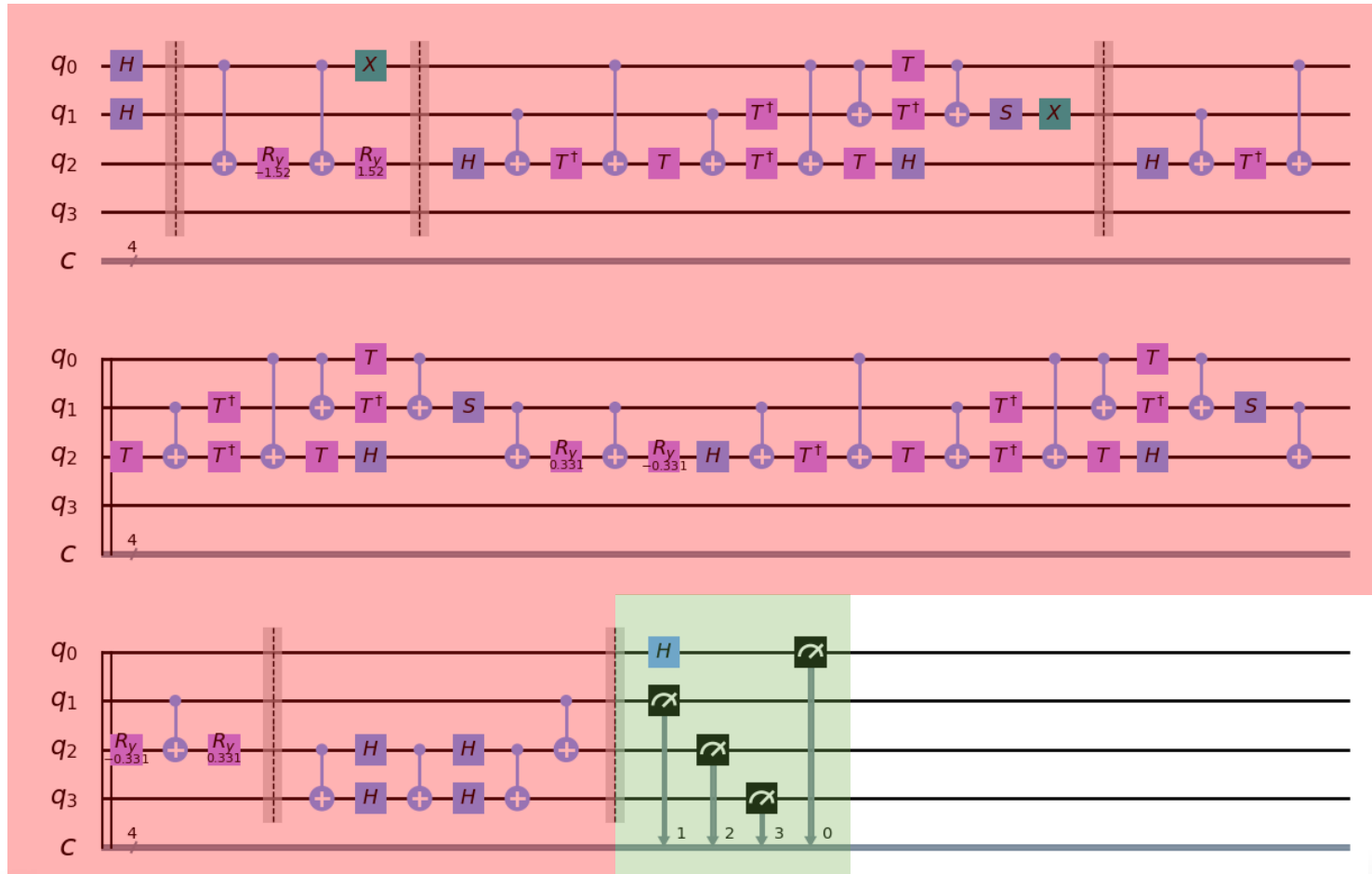
Aufbereiten der Daten

Klassifizierung



Kompakter Schaltkreis

Live Demo:





Support Vector Machine - Big Data (Ausblick)

Quantum Random Access Machine:

- Zugriff auf Superpositionen in $O(\log Nn)$
- Benötigt $O(nN)$ Hardware-Ressourcen

→ Q-SVM durchführbar in $O(\log Nn)$

SVM als „least squares“ Problem:

- Führe andere Lagrange-Multiplikatoren ein


$$c^p (\mathbf{w}^T \mathbf{v}^p + b) \geq 1 \rightarrow (\mathbf{w}^T \mathbf{v}^p + b) = c^p - c^p \beta_p$$

$$\Rightarrow \begin{bmatrix} 0 & \mathbf{1}^T \\ \mathbf{1} & \mathbf{K} + \gamma^{-1} \mathbf{I} \end{bmatrix} \begin{bmatrix} b \\ \boldsymbol{\alpha} \end{bmatrix} = \begin{bmatrix} 0 \\ \mathbf{c} \end{bmatrix}$$

- Löse Gleichungssystem mit HHL
- In $O(\text{polylog } Nn)$ durchführbar

Fazit:

- Größtes Bottleneck ist das Laden der Daten
- Durch Kernel-Trick lässt sich klassisch bereits in hohen Dimensionen arbeiten
- Datenkomprimierung ermöglicht aber Vorteil
- Bereits gute Ergebnisse auf NISQ-Rechnern



Danke für die Aufmerksamkeit
→ Fragen und Diskussion