

Clusteranalyse – Evaluierung

Teil 1: Intrinsische Evaluierung

Bewertung von Ergebnissen beim Clustern

Verschiedene Initialisierungen können unterschiedliche Ergebnisse liefern.
Um ein „bestes“ Ergebnis auszuwählen, braucht man ein Evaluationsmaß.



„Intrinsische Evaluierung“

d.h. es werden keine äußeren Faktoren berücksichtigt

- Clusteranalyse gehört zu den nicht überwachten Lernverfahren
- D.h. es gibt kein „korrektes“ Ergebnis, mit dem man das Ergebnis vergleichen kann

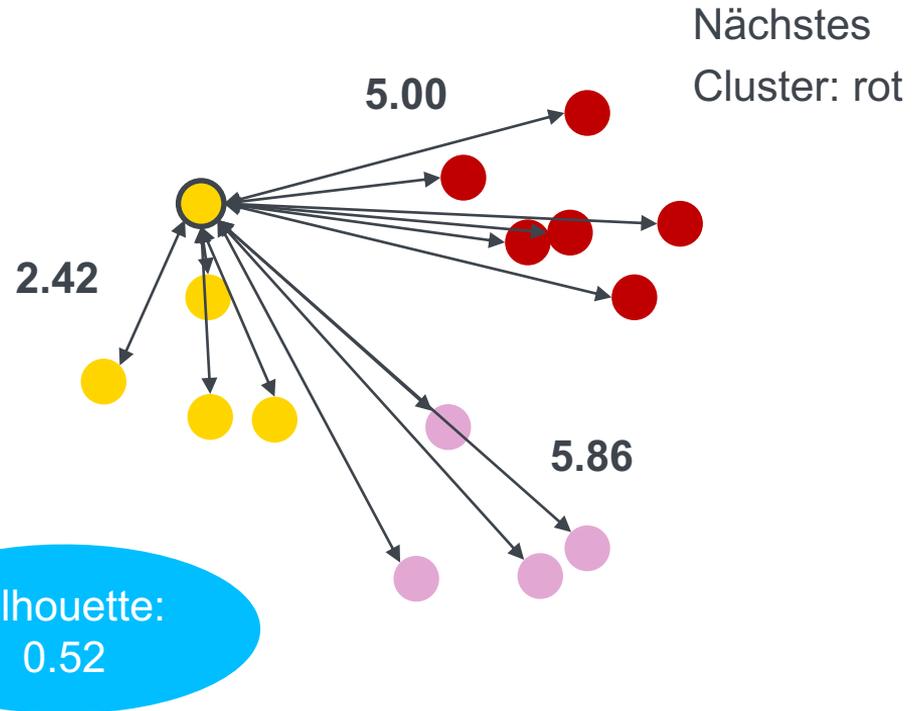
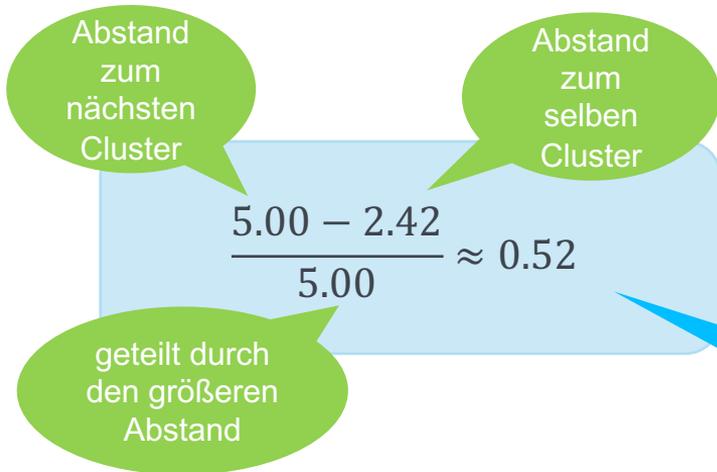
Bewerte, wie sauber die Cluster voneinander getrennt sind

Sind alle Punkte im passendsten Cluster?

Intrinsische Evaluierung

Variante 1: ein „gute“ Cluster-Struktur

- Idee:
 - Abstand zum selben Cluster: klein
 - Abstand zum nächsten Cluster: groß



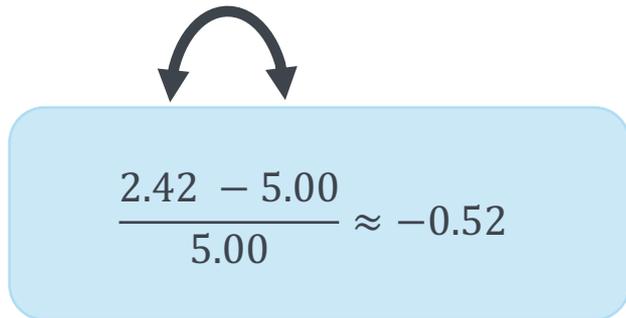
**Die Silhouette
(auch: Silhouettenkoeffizient)
eines Datenpunkts erfasst, wie gut
er in sein Cluster passt.**

**Ein Wert von 1 drückt aus,
dass der Datenpunkt perfekt in sein
Cluster passt.**

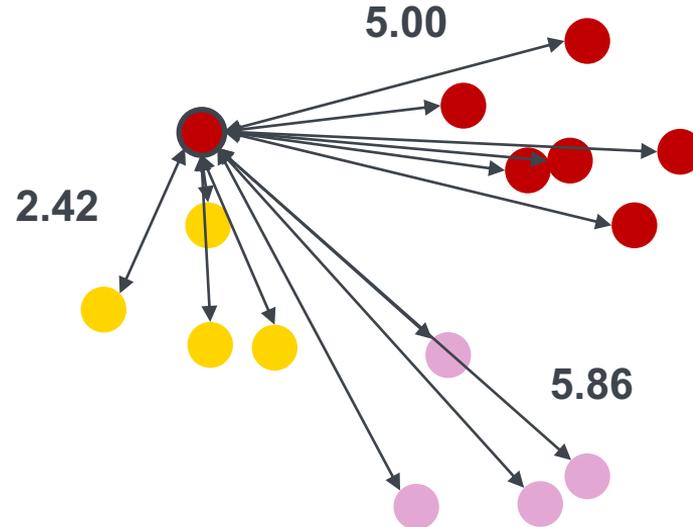
Man berechnet die Silhouette, indem man den mittleren Abstand zu Punkten des eigenen Clusters vom mittleren Abstand zu Punkten des nächsten Clusters abzieht und durch den größeren der beiden Werte teilt.

Variante 2: Ein „schlechte“ Cluster-Struktur

- Die Silhouette ist abhängig von der Zuordnung zum Cluster



$$\frac{2.42 - 5.00}{5.00} \approx -0.52$$

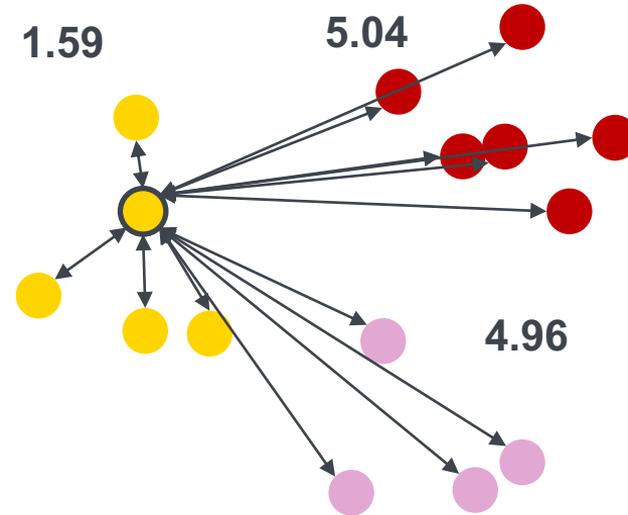


Ein negativer Wert für die Silhouette bedeutet, dass der Punkt näher zu Punkten eines anderen Clusters liegt als zu den Punkten des eigenen Clusters.

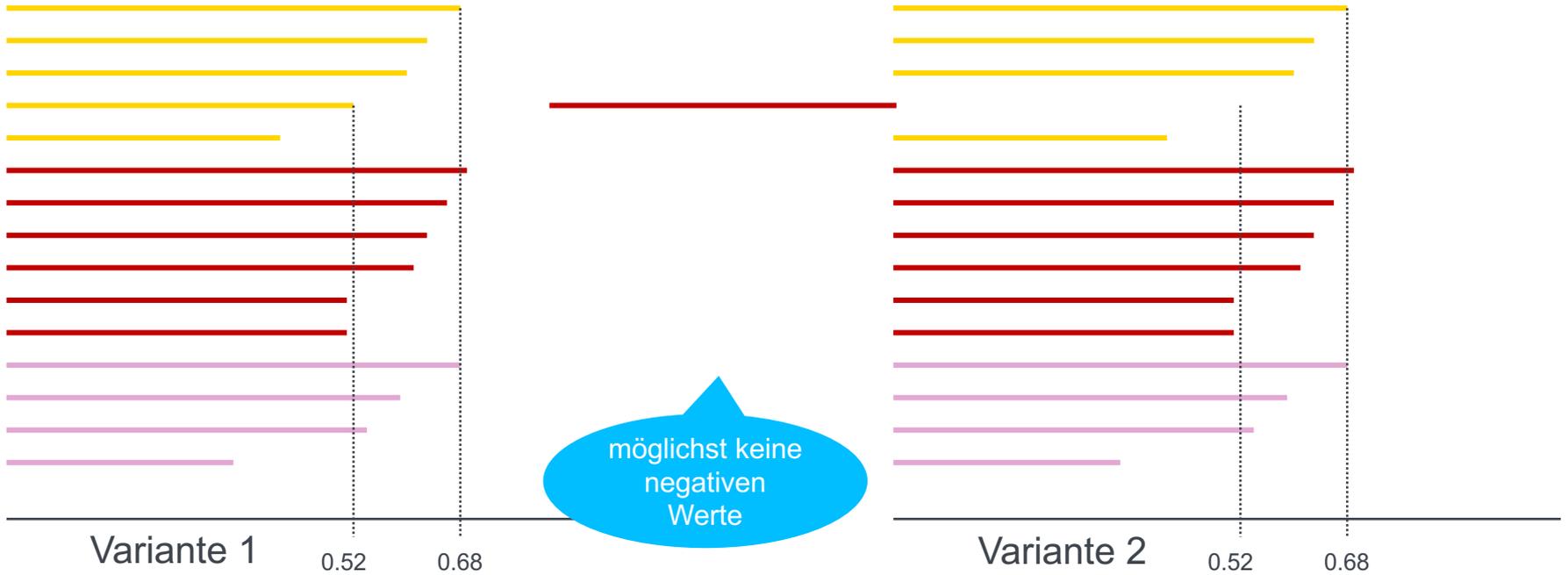
Wann ist ein Punkt im „richtigen“ Cluster?

- Silhouette für einen Punkt mitten im Cluster

$$\frac{4.96 - 1.59}{4.96} \approx 0.68$$



Der Silhouettenplot



möglichst keine
negativen
Werte

möglichst große Werte

Der Silhouettenplot zeigt die Silhouetten aller Datenpunkte an, sortiert nach Clustern. Innerhalb der Cluster werden die Werte nach der Größe sortiert.

Clusterergebnisse sind gut, wenn alle Werte möglichst groß sind und (fast) keine negativen Werte vorkommen.

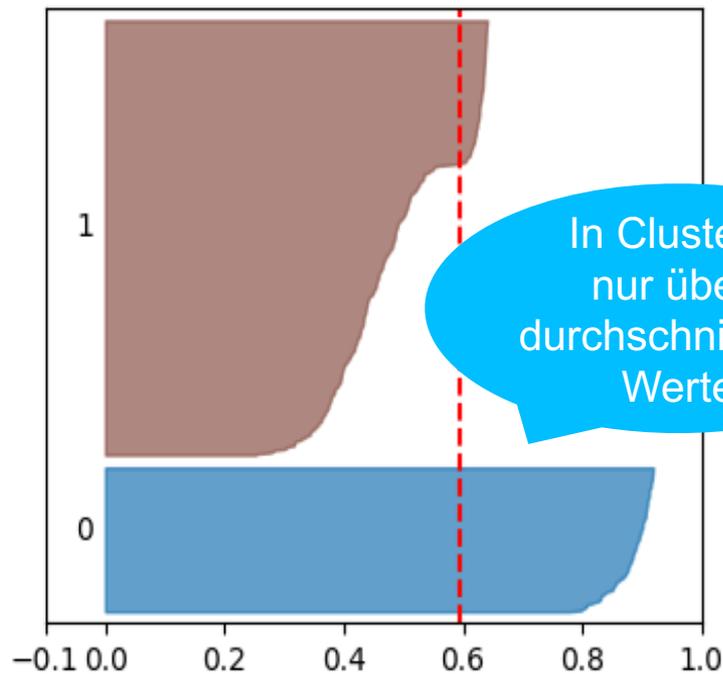
Ein Experiment

- 4 Punktwolken zufällig erzeugen lassen
- Mit k-Means mit verschiedenen Werten für k clustern lassen
- Silhouetten-Plots und Mittelwerte vergleichen

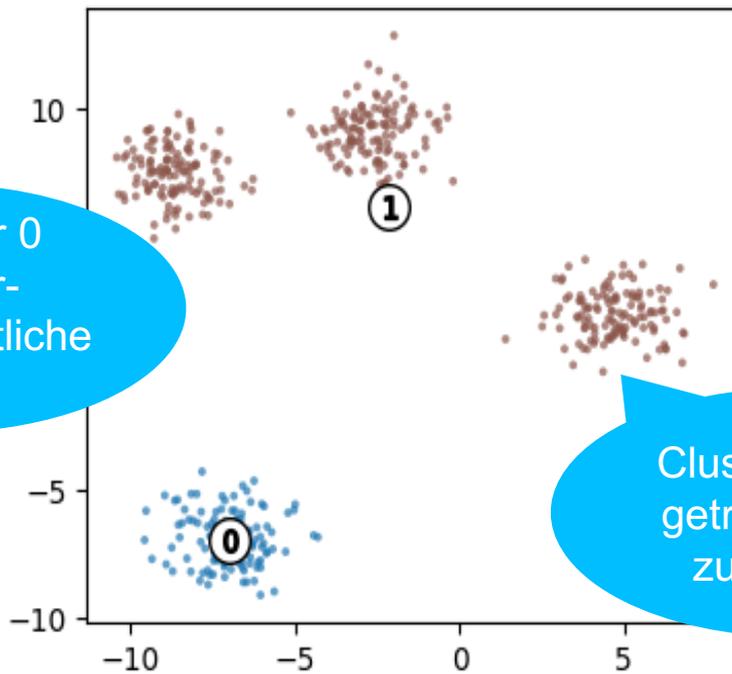


Experiment 1, k=2

Keine negativen Werte,
Mittelwert 0.60,
Silhouette von Cluster 1
schlechter



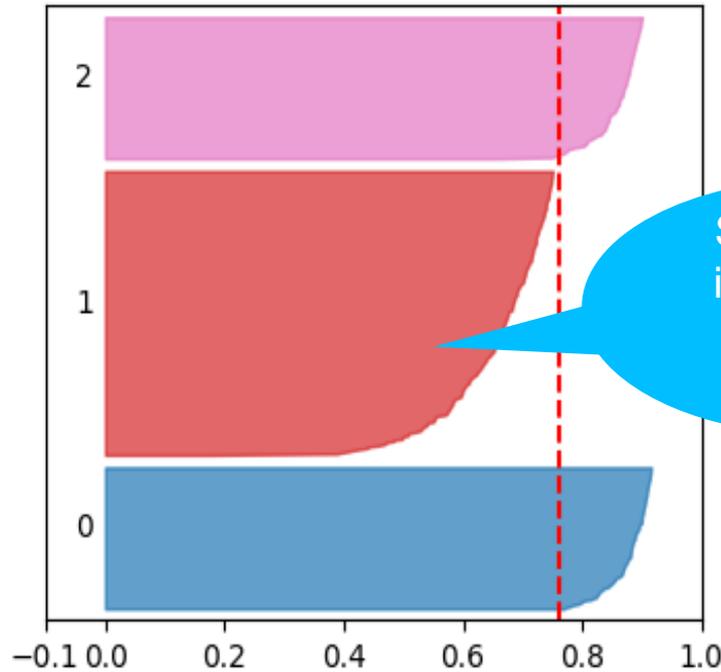
In Cluster 0
nur über-
durchschnittliche
Werte



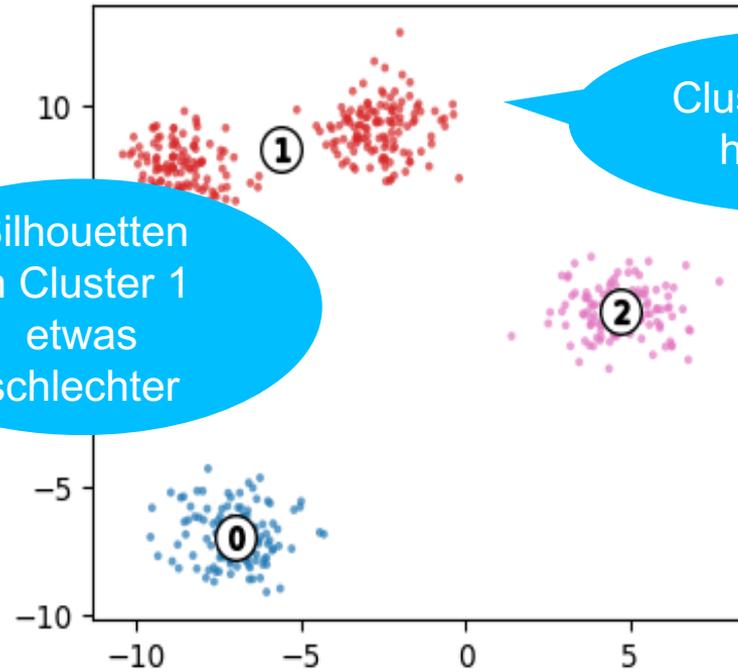
Cluster sauber
getrennt (aber
zu wenige)

Experiment 1, k=3

Mittelwert größer: 0.76



Silhouetten
in Cluster 1
etwas
schlechter

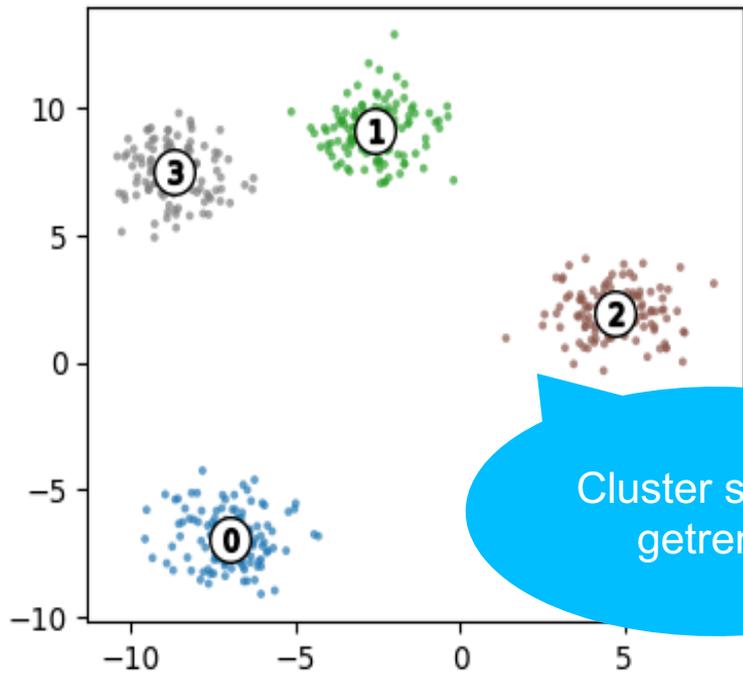
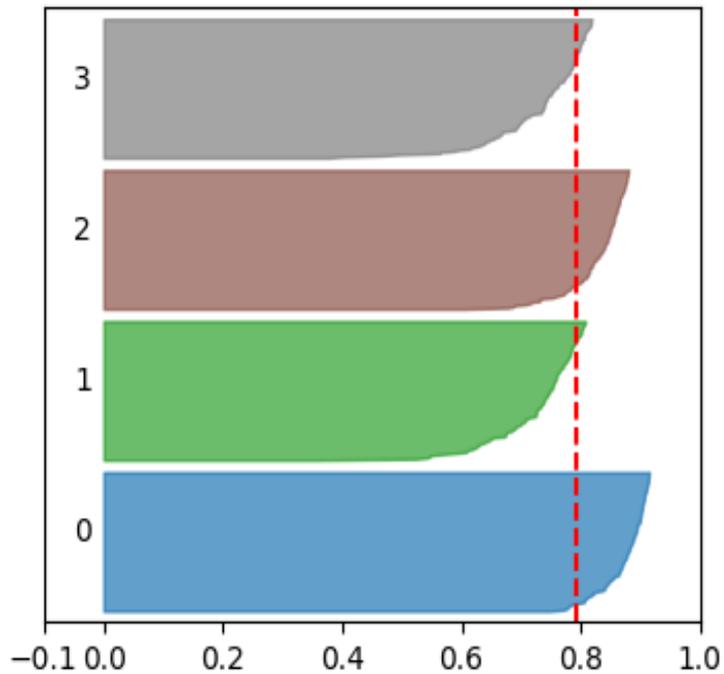


Cluster 1 nicht
homogen

Experiment 1, k=4

In allen Clustern
kommen ähnlich
gute Werte vor

Mittelwert noch größer:
0.79

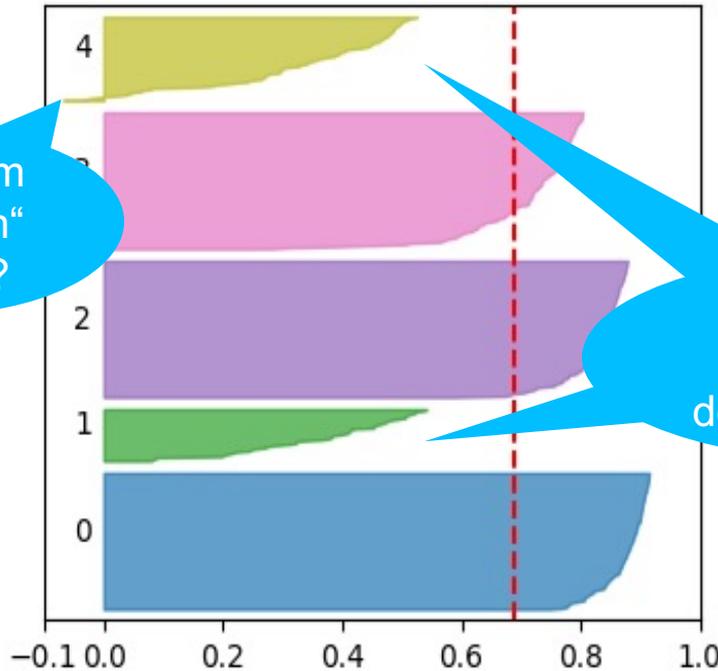


Cluster sauber
getrennt

Experiment 1, k=5

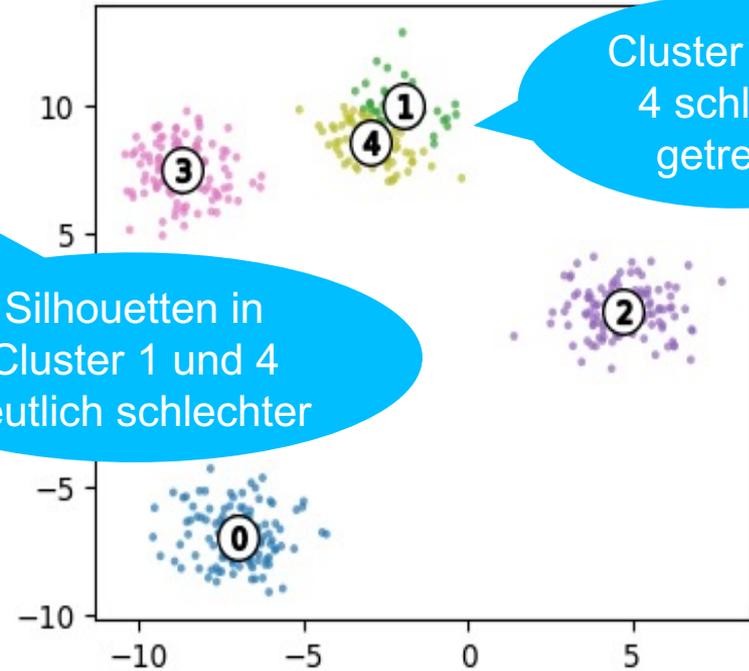
Silhouetten kleiner
(Mittelwert: 0.69),
einige negative Werte

Punkte im
„falschen“
Cluster?



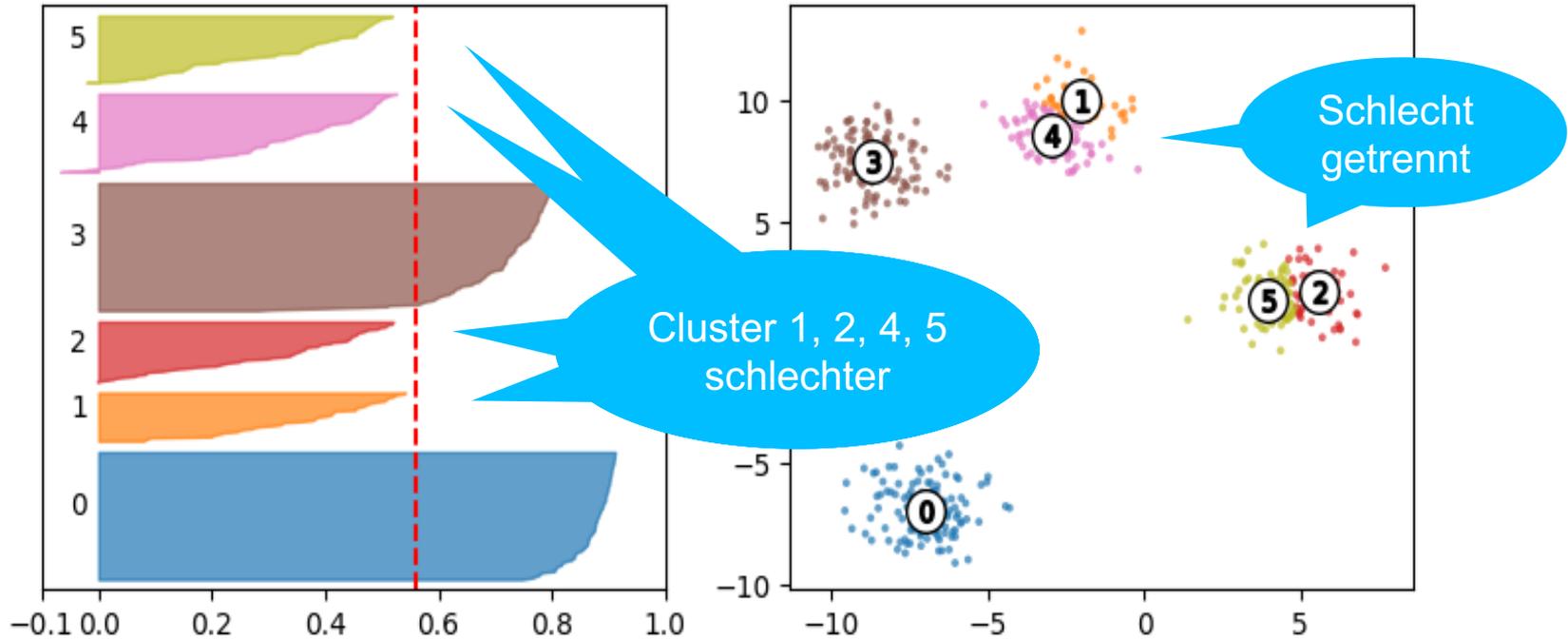
Silhouetten in
Cluster 1 und 4
deutlich schlechter

Cluster 1 und
4 schlecht
getrennt



Silhouetten noch kleiner
(Mittelwert: 0.56)

Experiment 1, k=6



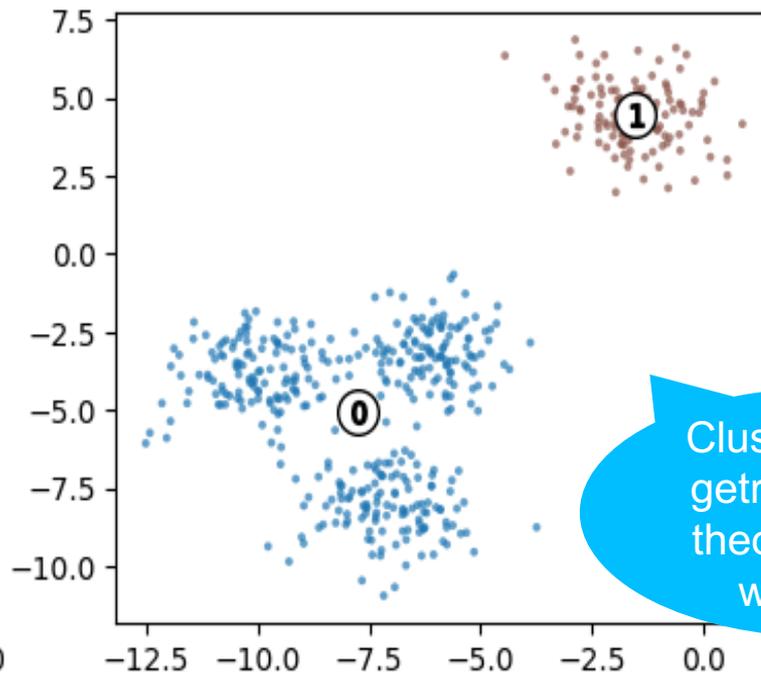
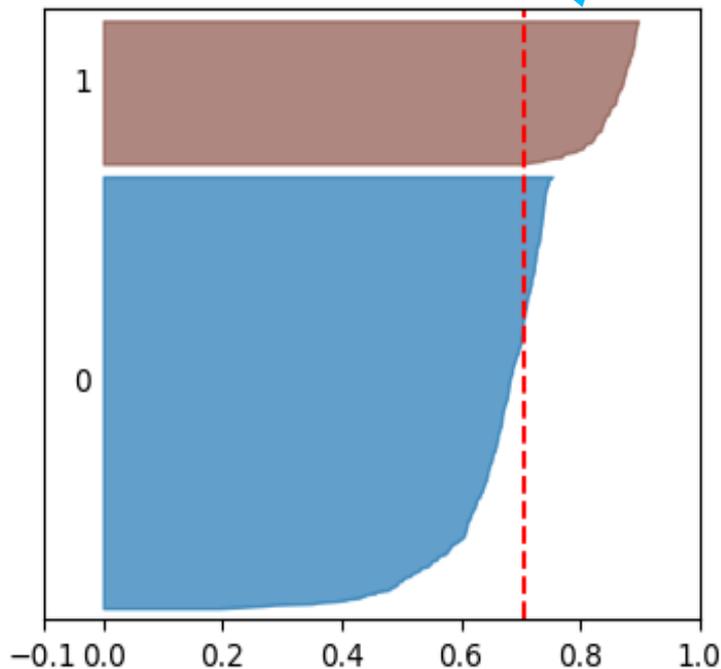
Zusammenfassung Experiment 1

- Bei $k=4$
 - höchster Mittelwert
 - alle Cluster ähnlich „gut“
- Bei $k=2$ und $k=3$
 - mäßige Unterschiede zwischen den Silhouetten der Cluster
- Bei $k=5$ und $k=6$
 - negative Werte, deutlich niedrigere Mittelwerte
 - große Unterschiede zwischen den Silhouetten der Cluster

Experiment 2, k=2

In Cluster 1
fast nur über-
durchschnittliche
Werte

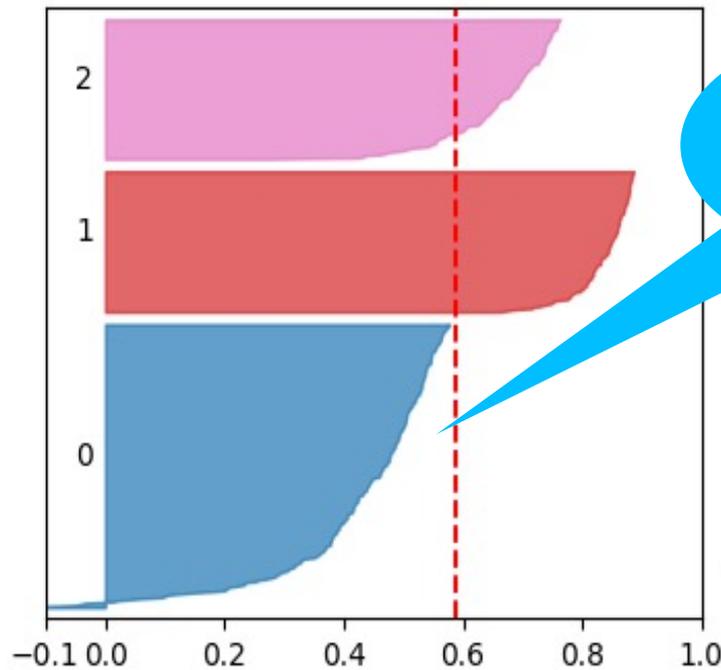
Keine negativen Werte,
Silhouetten relativ groß
(Mittelwert: 0.70)



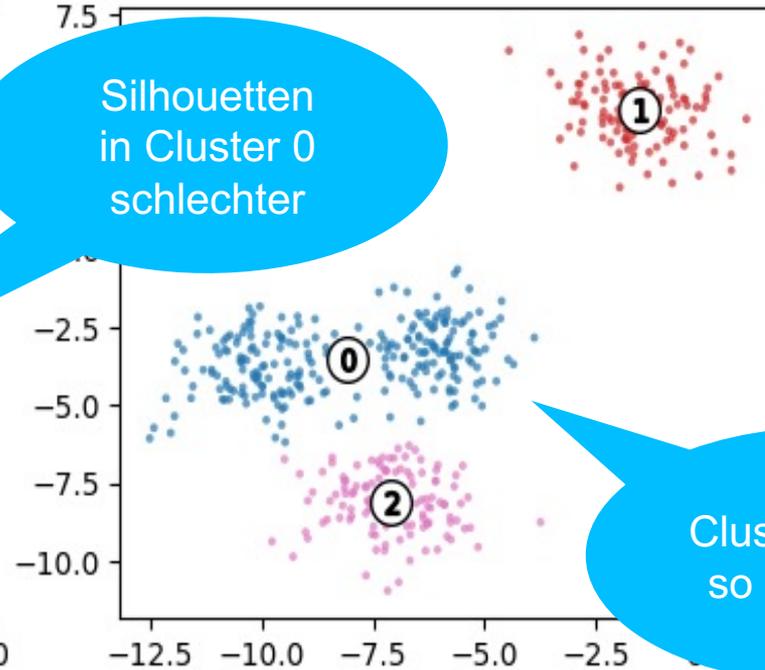
Cluster sauber
getrennt (aber
theoretisch zu
wenige!!)

Experiment 2, k=3

Negative Werte,
Silhouetten insgesamt
kleiner (Mittelwert: 0.59)

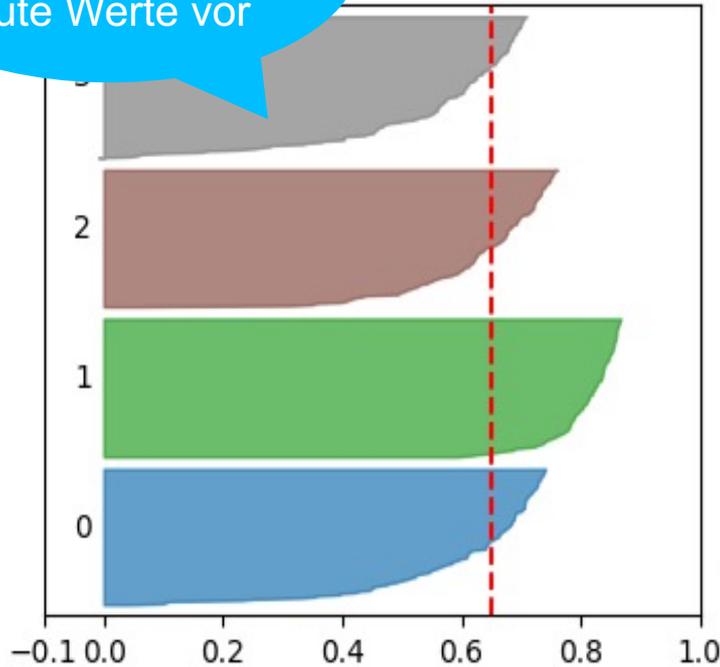


Silhouetten
in Cluster 0
schlechter

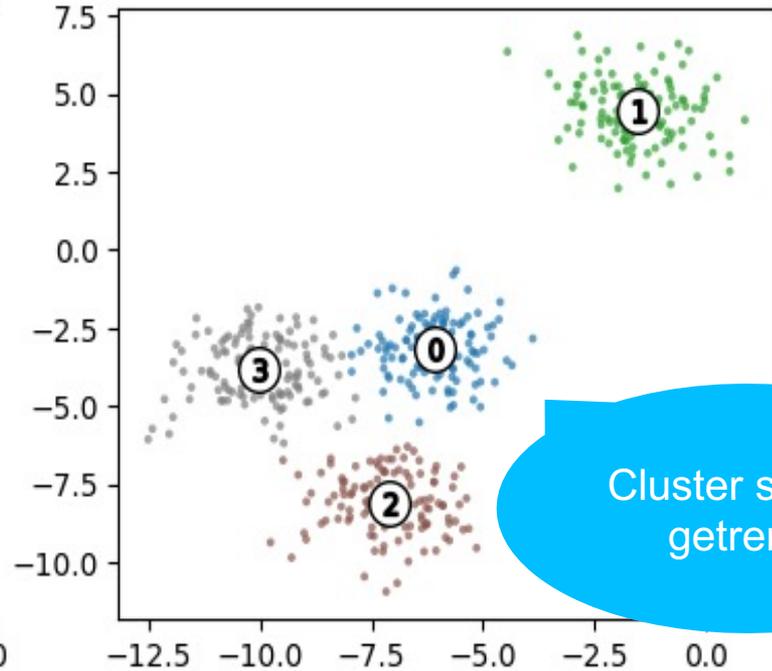


Cluster 0 nicht
so homogen

In allen Clustern
kommen ähnlich
gute Werte vor



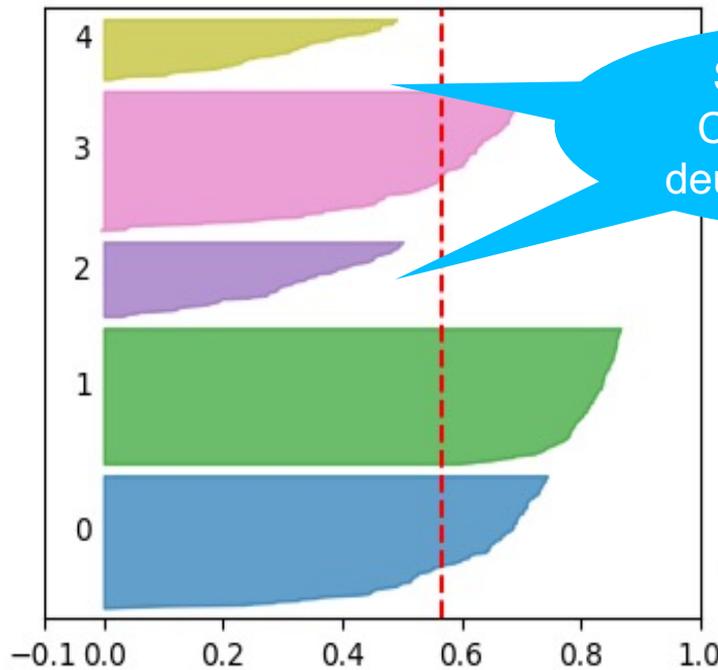
Keine negativen Werte,
Silhouetten relativ groß
(Mittelwert: 0.65)



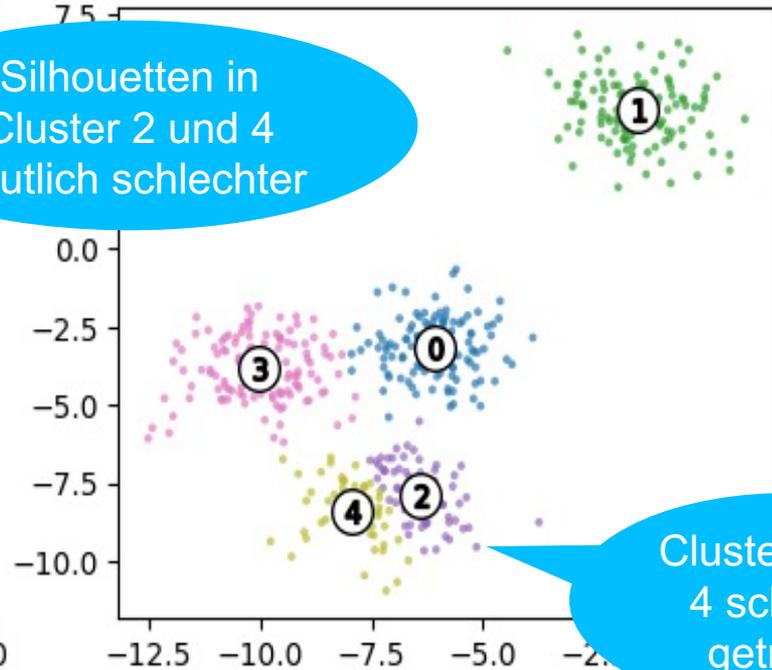
Cluster sauber
getrennt

Experiment 2, k=5

Keine negativen Werte,
Silhouetten kleiner
(Mittelwert: 0.56)



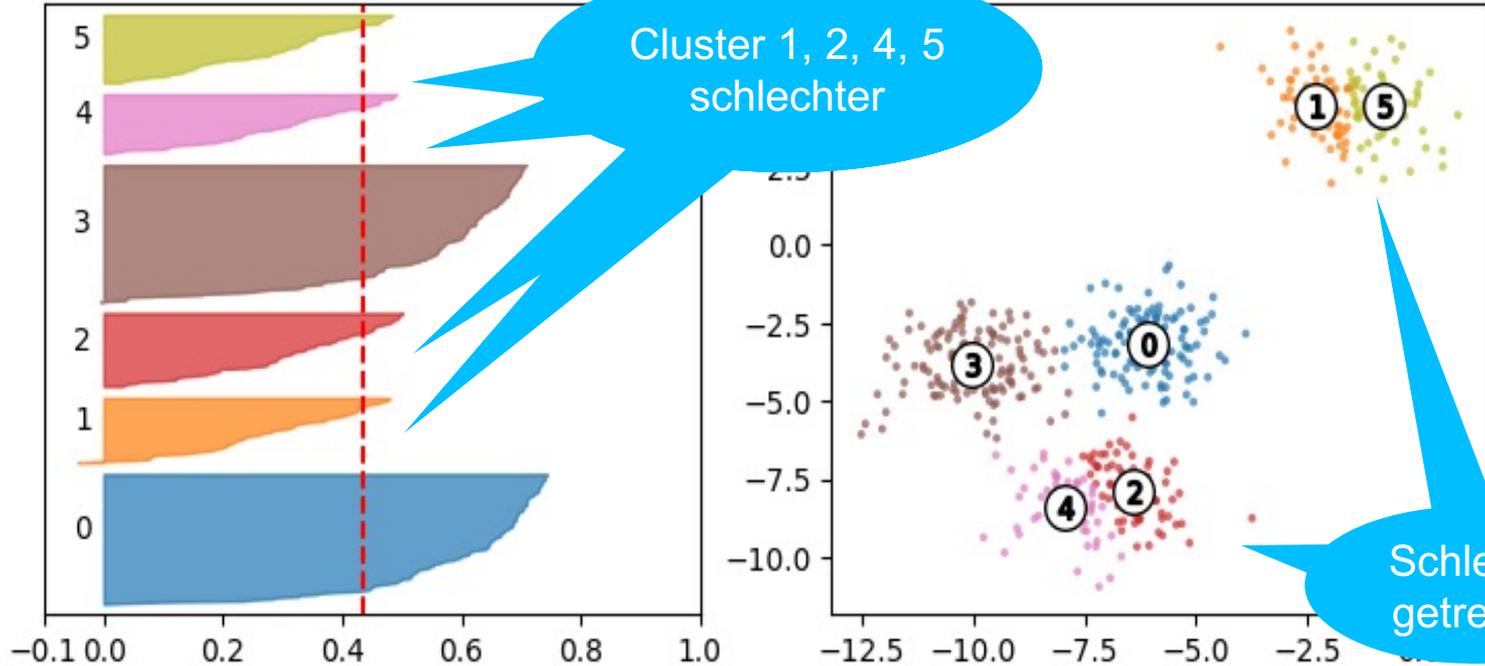
Silhouetten in
Cluster 2 und 4
deutlich schlechter



Cluster 2 und
4 schlecht
getrennt

Experiment 2, k=6

Silhouetten noch kleiner
(Mittelwert: 0.48)



Zusammenfassung Experiment 2

- Unerwartet: bei $k=2$ höchster Mittelwert!!!
 - vielleicht sinnvoll, da drei der vier Cluster nahe beieinander
- Immerhin: sowohl bei $k=2$ als auch bei $k=4$ keine negativen Werte
- Bei $k=4$ ähnlichste Silhouetten für alle Cluster

Der Silhouetten-Mittelwert sowie Silhouettenplots dienen der intrinsischen Evaluierung von Clusterergebnissen.

Dabei sollte ein insgesamt hoher Mittelwert angestrebt werden. Außerdem ist es sinnvoll, Cluster zu vermeiden, die deutlich schlechtere Silhouetten aufweisen.

Silhouettenanalyse

- Kann zur intrinsischen Evaluierung verwendet werden
- Es gibt weitere Maße zu diesem Zweck!
- Silhouetten sind wahrscheinlich das bekannteste Maß

Dr. Antje Schweitzer

Universität Stuttgart
Institut für Maschinelle Sprachverarbeitung



Universität Stuttgart

Institut für Maschinelle Sprachverarbeitung
Institut für Software Engineering



Industrie- und Handelskammer
Reutlingen

Reutlingen | Tübingen | Zollernalb



Region Stuttgart



Industrie- und Handelskammer
Karlsruhe



Lizenzbestimmungen

“Clusteranalyse – Evaluierung, Teil 1: Intrinsische Evaluierung“ von Antje Schweitzer, KI B³ / Uni Stuttgart

Das Werk - mit Ausnahme der folgenden Elemente:

- Logos der Verbundpartner und des Förderprogramms
- im Quellenverzeichnis aufgeführte Medien

ist lizenziert unter:

 [CC BY 4.0 \(https://creativecommons.org/licenses/by/4.0/deed.de\)](https://creativecommons.org/licenses/by/4.0/deed.de)

(Namensnennung 4.0 International)

Quellenverzeichnis

Titelfoto: [Omar Flores \(https://unsplash.com/de/@designedbyflores\)](https://unsplash.com/de/@designedbyflores), ohne Titel, auf [Unsplash \(https://unsplash.com/de/fotos/blaue-rote-und-weisse-grafiken-IQT_bOWtysE\)](https://unsplash.com/de/fotos/blaue-rote-und-weisse-grafiken-IQT_bOWtysE), lizenziert unter [Unsplash-Lizenz \(https://unsplash.com/license\)](https://unsplash.com/license). Bildausschnitt verändert.