

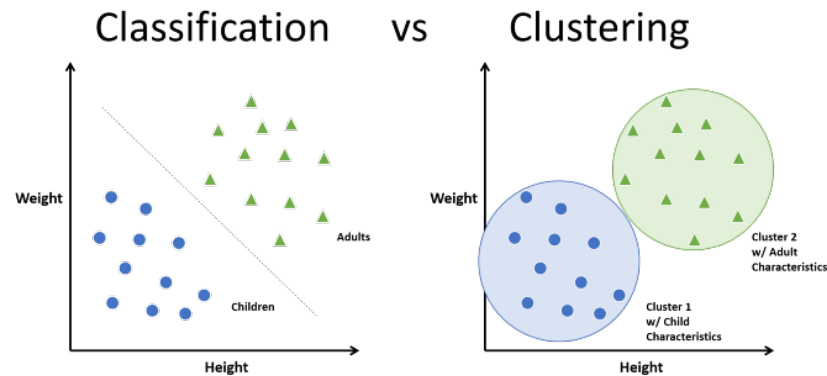


VL 3 (Teil1) – Klassifikation / Clusteranalyse

Mareike Schulze

Klassifikation vs. Clusteranalyse

- Klassifikation:
 - Zuordnung von Datenpunkten zu vorhandenen Klassen
 - anhand von Ähnlichkeiten bzw. übereinstimmender Merkmale
 - supervised learning

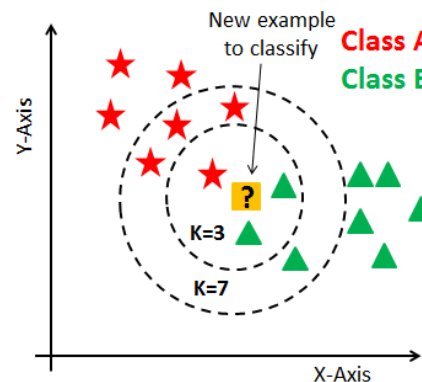


- Clusteranalyse:
 - Identifizierung von Ähnlichkeiten unter Datenstrukturen
 - Gruppierung dieser ähnlichen Datenpunkte
 - unsupervised learning

k-Nearest Neighbors (kNN) – Klassifikation

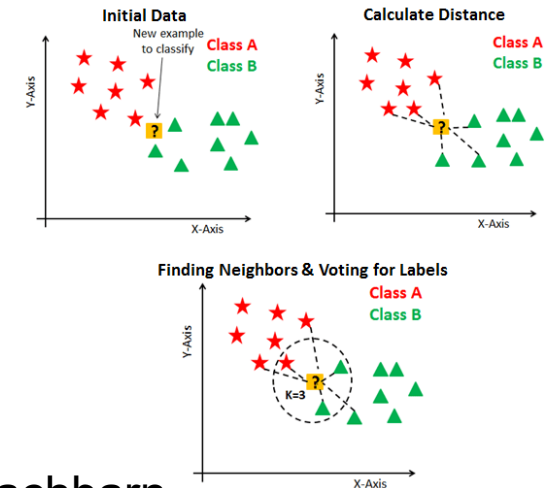
- Vorgehen:

1. Berechnet die Distanz eines Punktes zu den Nachbarn
2. Identifiziert die k Nachbarn mit den geringsten Abständen
3. Entscheidet anhand der k Nachbarn über die Zugehörigkeit des betrachteten Punktes (welche Zugehörigkeit haben die Nachbarn?)



- Parameter (Auszug):

- k : Anzahl zu untersuchender Nachbarn
- Gewichtung: gleichwertig oder nach Distanz
- Algorithmus: BallTree, KDTree, brute-force, auto

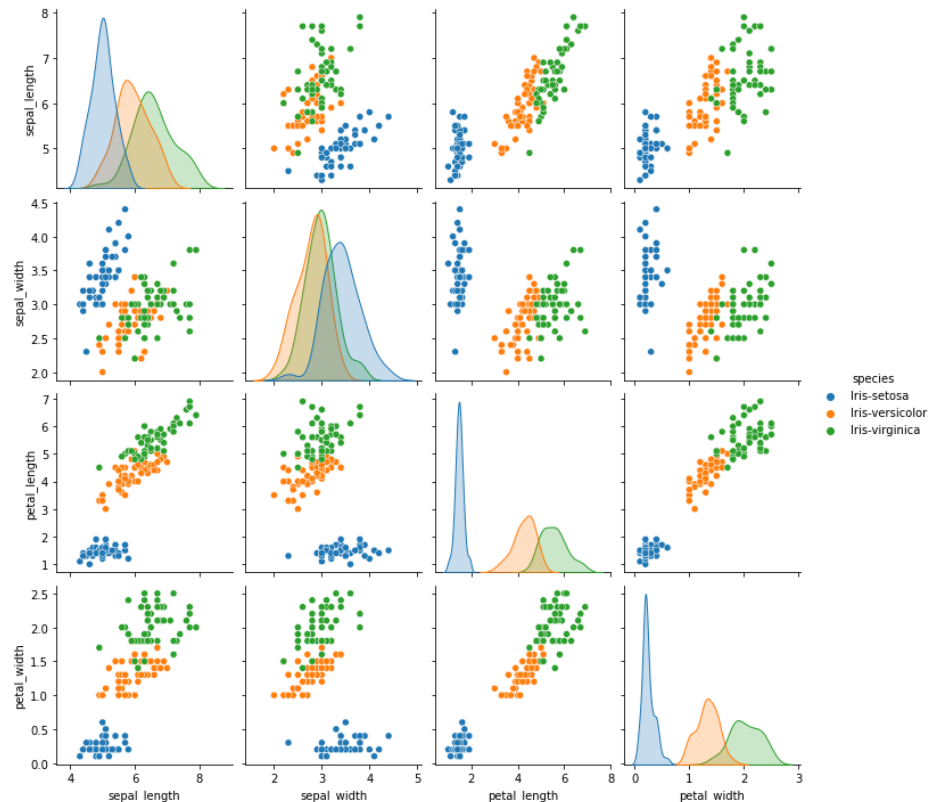


Exkurs: Eine erste Übersicht über den Datensatz

```
1 # Bibliotheken laden
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Datensatz aus der Bibliothek seaborn laden
6 df = sns.load_dataset("iris")
7
8 # Eine erste Übersicht über den Datensatz
9 print(df.head())
10 print(df.describe())
11
12 # Den Pairplot erstellen und anzeigen
13 sns.pairplot(df, hue="species")
14 plt.show()
```

describe()

	sepal_length	sepal_width	petal_length	petal_width
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.054000	3.758667	1.198667
std	0.828066	0.433594	1.764420	0.763161
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000



head()

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

k-Nearest Neighbors – Lernbeispiel [1/2]

```

1 from sklearn.datasets import load_iris
2 from sklearn.model_selection import train_test_split
3 from sklearn.neighbors import KNeighborsClassifier
4 from sklearn.metrics import confusion_matrix, accuracy_score
5 import knnplot
6
7 # Datensatz laden
8 iris = load_iris()
9
10 # Schlüssel anzeigen (Datenstruktur ist Dictionary)
11 print(iris.keys())
12
13 # Auswählen und Anzeigen der Features
14 X = iris.data
15 print(iris.feature_names)
16 print(iris.data.tolist())
17
18 # Auswählen und Anzeigen der Zielvariablen
19 y = iris.target
20 print(iris.target_names)
21 print(iris.target)
22
23 # Datensatz aufteilen in Trainings- und Testdaten
24 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=142)
25 print("X_train:", X_train.shape, "X_test:", X_test.shape)
26 print("y_train:", y_train.shape, "y_test:", y_test.shape)

```


k-Nearest Neighbors – Lernbeispiel [1/2]

```
28 # Den kNN-Klassifikator aussuchen und anwenden
29 kNN = KNeighborsClassifier(n_neighbors=5)
30 kNN.fit(X_train, y_train)
31
32 # Ausgabe der zu den einzelnen Datenpunkten zugeordneten Klassen
33 y_pred = kNN.predict(X_test)
34 print(y_pred)
35 # direkter Vergleich mit den tatsächlichen Klassen
36 print(y_test)
37
38 # Beispiel: Ausgabe der Klassenzugehörigkeit eines Datenpunktes
39 print(iris.target_names[y_pred[14]])
40 # direkter Vergleich mit der tatsächlichen Klasse
41 print(iris.target_names[y_test[14]])
42
43 # Beispiel: Ausgabe der Wahrscheinlichkeiten der Klassenzugehörigkeit
44 # hier: 60% virginica, 40% versicolor
45 print(kNN.predict_proba(X_test)[14])
46
47 # Erzeugen und Anzeigen der Konfusionsmatrix
48 print(confusion_matrix(y_test, y_pred))
49
50 # Die Genauigkeit ermitteln und ausgeben
51 accuracy = accuracy_score(y_test, y_pred)*100
52 print('Accuracy of our model is equal ' + str(round(accuracy, 2)) + ' %.')
```

array([0, 1, 1, 2, 1, 1, 0, 0, 2, 1, 1, 1, 2, 0, 2, 0, 2, 1, 1, 2, 2, 1,
0, 1, 2, 2, 2, 2, 0, 1, 2, 1, 2, 2, 1, 2, 1, 2])

array([0, 1, 1, 2, 1, 1, 0, 0, 2, 1, 1, 1, 2, 0, 1, 0, 2, 1, 2, 1, 2, 1,
0, 1, 2, 2, 2, 2, 0, 1, 2, 1, 2, 1, 1, 2, 1, 2])

'virginica' ←

'versicolor' ←

array([0. , 0.4, 0.6])

[[7 0 0]
[0 14 3]
[0 1 13]]

Alle *setosa* wurden korrekt klassifiziert.
3 *versicolor* wurden als *virginica* klassifiziert.
1 *virginica* wurde als *versicolor* klassifiziert.

'Accuracy of our model is equal 89.47 %.'

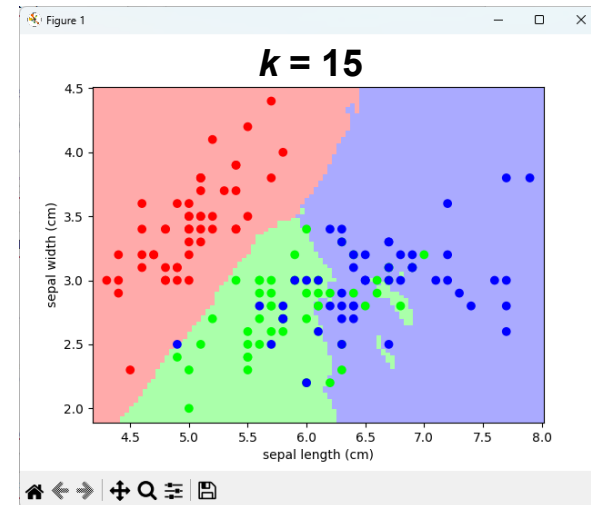
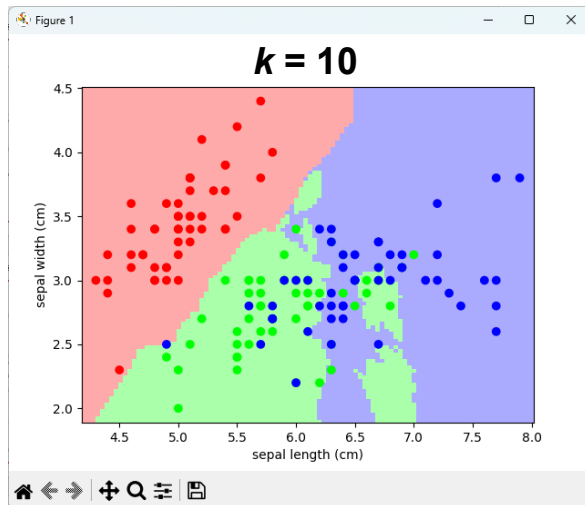
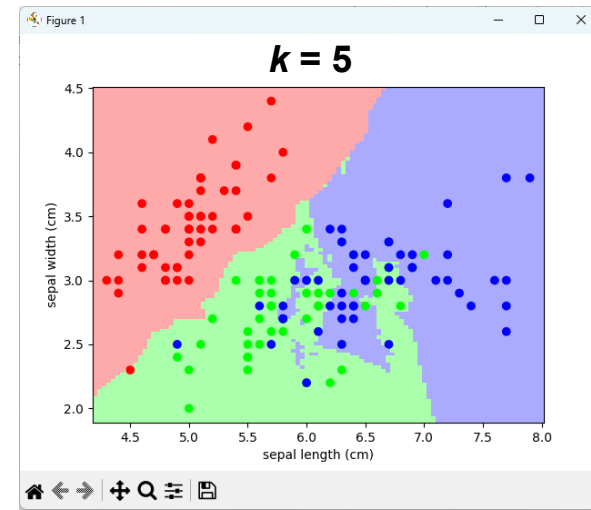
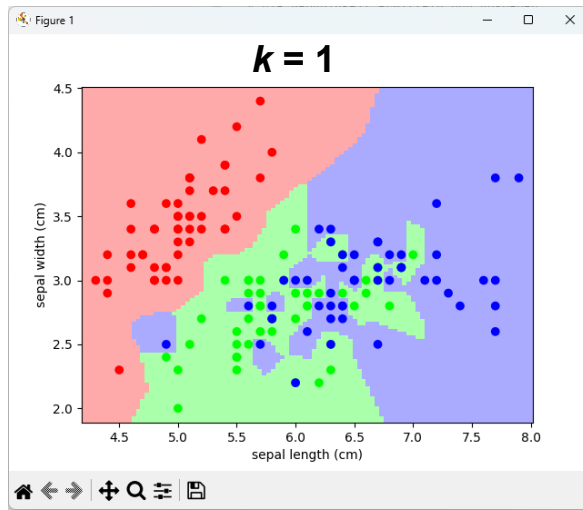
k-Nearest Neighbors – Abhängigkeit von k [1/2]

```
54 # Die Klassifikation über 2 Parameter (sepal_length und sepal_width)
55 # jeweils mit k = 1, 5, 10, 15
56 X = iris.data[:, :2]
57
58 knn = KNeighborsClassifier(n_neighbors=1)
59 knn.fit(X, y)
60 knnplot.plot_iris_knn(knn, X, y)
61
62 knn = KNeighborsClassifier(n_neighbors=5)
63 knn.fit(X, y)
64 knnplot.plot_iris_knn(knn, X, y)
65
66 knn = KNeighborsClassifier(n_neighbors=10)
67 knn.fit(X, y)
68 knnplot.plot_iris_knn(knn, X, y)
69
70 knn = KNeighborsClassifier(n_neighbors=15)
71 knn.fit(X, y)
72 knnplot.plot_iris_knn(knn, X, y)
```

Inhalt der Datei *kNNplot.py*

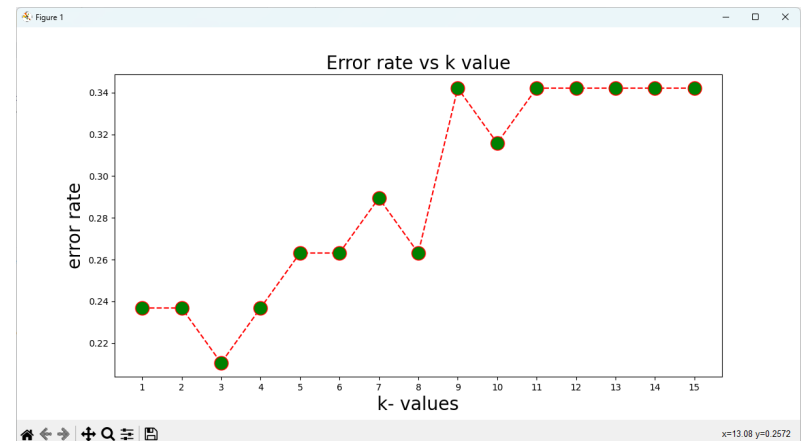
```
1 import numpy as np
2 import pylab as pl
3 import matplotlib.pyplot as plt
4 from matplotlib.colors import ListedColormap
5
6 cmap_light = ListedColormap(['#FFAAAA', '#AAFFAA', '#AAAAFF'])
7 cmap_bold = ListedColormap(['#FF0000', '#00FF00', '#0000FF'])
8
9
10 4 usages
11 def plot_iris_knn(knn, X, y):
12     x_min, x_max = X[:, 0].min() - .1, X[:, 0].max() + .1
13     y_min, y_max = X[:, 1].min() - .1, X[:, 1].max() + .1
14     xx, yy = np.meshgrid(*xi: np.linspace(x_min, x_max, num: 100),
15                           np.linspace(y_min, y_max, num: 100))
16     Z = knn.predict(np.c_[xx.ravel(), yy.ravel()])
17
18     # Visualisierung der Klassifikation (Klassengrenzen)
19     Z = Z.reshape(xx.shape)
20     pl.figure()
21     pl.pcolormesh(*args: xx, yy, Z, cmap=cmap_light)
22
23     # Einzeichnen der einzelnen Datenpunkte (tatsächliche Klassen)
24     pl.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap_bold)
25     pl.xlabel('sepal length (cm)')
26     pl.ylabel('sepal width (cm)')
27     pl.axis('tight')
28
29     plt.show()
```

k-Nearest Neighbors – Abhängigkeit von k [2/2]



k-Nearest Neighbors – Das optimale k

```
1 from sklearn.datasets import load_iris
2 from sklearn.model_selection import train_test_split
3 from sklearn.neighbors import KNeighborsClassifier
4 import numpy as np
5 import matplotlib.pyplot as plt
6
7 iris = load_iris()
8
9 # zur Veranschaulichung wieder nur 2 Features (sepal_length und sepal_width)
10 X = iris.data[:, :2]
11 y = iris.target
12
13 X_train, X_test, y_train, y_test = train_test_split(*arrays: X, y, random_state=142)
14
15 # Das optimale k ermitteln
16 error_rate = []
17 kMin, kMax = 1, 15
18
19 for i in range(kMin, kMax+1):
20     KNN = KNeighborsClassifier(n_neighbors=i)
21     KNN.fit(X_train, y_train)
22     y_pred = KNN.predict(X_test)
23     error_rate.append(np.mean(y_pred != y_test))
24
25 # Die Ergebnisse ausgeben
26 plt.figure(figsize=(12, 6))
27 plt.plot(*args: range(kMin, kMax+1), error_rate, marker="o", markerfacecolor="green",
28         linestyle="dashed", color="red", markersize=15)
29 plt.title(label: "Error rate vs k value", fontsize=20)
30 plt.xlabel(xlabel: "k- values", fontsize=20)
31 plt.ylabel(ylabel: "error rate", fontsize=20)
32 plt.xticks(range(kMin, kMax+1))
33 plt.show()
```



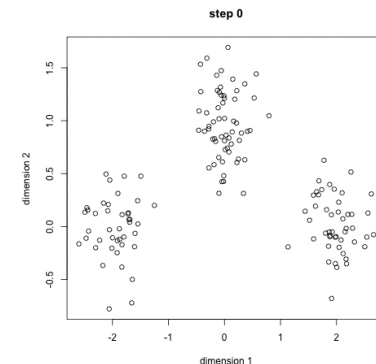
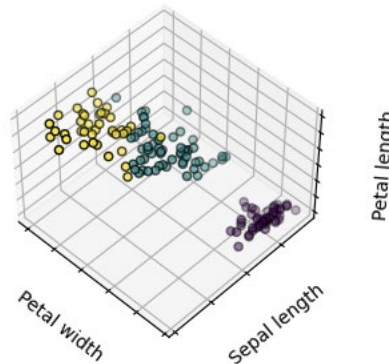
k-Nearest Neighbors – Fazit

- Auch als Regression möglich
- Vorteile:
 - Einfache Anwendung (kann einfach erlernt/nachvollzogen werden)
 - Einfache Anpassung an neue Daten (nach demselben Muster)
 - Wenige notwendige Parameter (alles außer Daten und k ist optional)
- Nachteile:
 - Vergleichsweise langsam und ressourcenhungrig
 - Für hochdimensionale Daten nur schlecht geeignet
 - Anfälligkeit (niedriges k führt zu Overfitting,
hohes k führt zu Glättung / Unteranpassung)

k-Means / k-Median – Clusteranalyse

- Vorgehen:

1. Bestimme die Anzahl k der gewünschten Cluster
2. Bestimme die Mittelpunkte / Mediane der einzelnen Cluster (random)
3. Ordne jeden Datenpunkt dem nächstgelegenen Mittelpunkt / Median zu
4. Berechne die Mittelpunkte / Mediane der einzelnen Cluster neu
5. Wiederhole 3.-4., solange die Mittelpunkte / Mediane sich verändern



- Parameter (Auszug):

- k : Anzahl erwünschter Cluster
- n_init : Anzahl der Initialisierungen
- max_iter : Anzahl der Iterationen pro Initialisierung

k-Means – Lernbeispiel [1/2]

```
1 from sklearn.datasets import load_iris
2 from sklearn.cluster import KMeans
3 import kmplot
4
5 # Datensatz laden
6 iris = load_iris()
7
8 # Auswählen und Anzeigen der Features
9 X = iris.data
10 print(iris.feature_names) ['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']
11 print(iris.data.tolist()) [[5.1, 3.5, 1.4, 0.2], [4.9, 3.0, 1.4, 0.2], [4.7, 3.2, 1.3, 0.2], [4. ...]
12
13 # Den kNN-Klassifikator aussuchen und anwenden
14 kM = KMeans(n_clusters=3)
15 kM.fit(X)
16
17 # Die Cluster ausgeben
18 print(kM.labels_) # kM.predict(X)
19
20 # die Clusterzentren ausgeben
21 print(kM.cluster_centers_)
```

```
array([[1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
        2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2,
        2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 0, 2, 0, 0, 0, 0, 2, 0, 0,
        0, 0, 0, 2, 2, 0, 0, 0, 0, 2, 0, 2, 0, 2, 0, 0, 2, 2, 0, 0,
        0, 2, 0, 0, 0, 0, 2, 0, 0, 0, 2, 0, 0, 0, 2, 0, 0, 2])

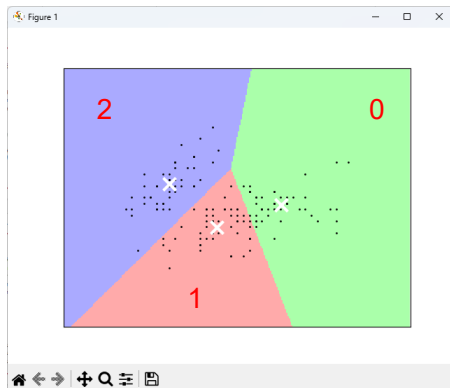
array([[6.85      , 3.07368421, 5.74210526, 2.07105263],
       [5.006      , 3.428      , 1.462      , 0.246      ],
       [5.9016129 , 2.7483871 , 4.39354839, 1.43387097]])
```

k-Means – Lernbeispiel [2/2]

```

23 # Das Clustering über 2 Parameter (sepal_length und sepal_width)
24 # jeweils mit k = 1, 5, 10, 15
25 X = iris.data[:, :2]
26
27 km = KMeans(n_clusters=3)
28 km.fit(X)
29
30 print(km.labels_)
31 print(km.cluster_centers_)
32
33 kmplot.plot_iris_km(km, X)

```

[illegible]

Inhalt der Datei *kMplot.py*

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from matplotlib.colors import ListedColormap
4
5 cmap_light = ListedColormap(['#FFAAAA', '#AAFFAA', '#AAAAFF'])
6
7
8 1 usage
9
10 def plot_iris_kM(kM, X):
11     x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
12     y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
13     xx, yy = np.meshgrid(*xi: np.arange(x_min, x_max, 0.02), np.arange(y_min, y_max, 0.02))
14     Z = kM.predict(np.c_[xx.ravel(), yy.ravel()])
15
16     # Visualisierung des Clustering (Clustergrenzen)
17     Z = Z.reshape(xx.shape)
18     plt.figure(1)
19     plt.clf()
20     plt.imshow(Z, interpolation="nearest", extent=(xx.min(), xx.max(), yy.min(), yy.max()),
21               cmap=cmap_light, aspect="auto", origin="lower")
22     plt.plot(*args: X[:, 0], X[:, 1], "k.", markersize=2)
23
24     # Visualisierung der Clusterzentren
25     centroids = kM.cluster_centers_
26     plt.scatter(centroids[:, 0], centroids[:, 1], marker="x", s=169, linewidths=3,
27               color="w", zorder=10)
28     plt.xlim(*args: x_min, x_max)
29     plt.ylim(*args: y_min, y_max)
30     plt.xticks(())
31     plt.yticks(())
32     plt.show()

```

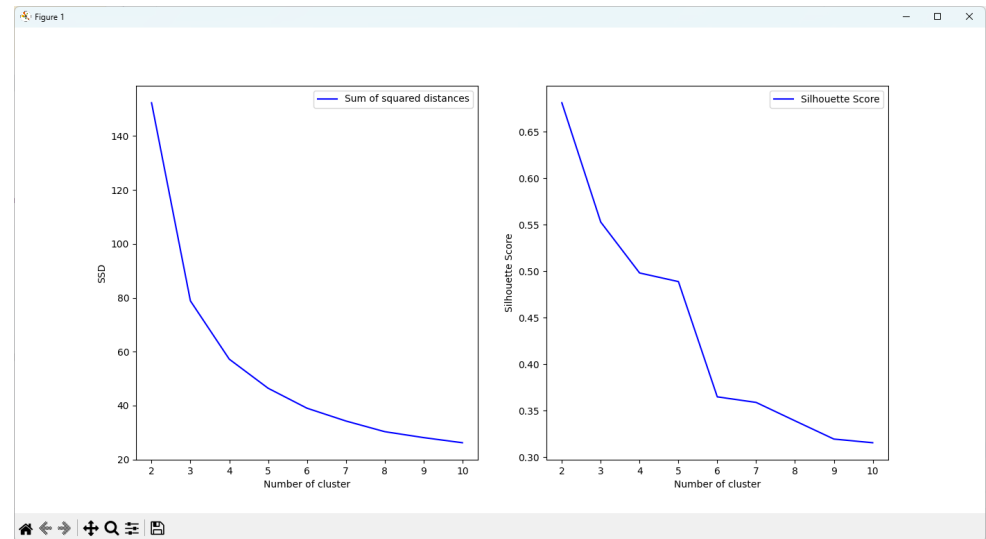
kMeans – Das optimale k

- Mittels Ellenbogenkurve
 - hier wird ermittelt, ab wann sich ein höherer Aufwand (ein höheres k) nicht mehr rechnet bzw. das Clustering nicht mehr verbessert
 - es wird pro k die Summe der Abstände der einzelnen Datenpunkte zu Ihrem nächsten Clusterzentrum zum Quadrat berechnet
 - dadurch kann abgelesen werden, welche Vergrößerung des k welche Verbesserung dieses Maßes hervorruft
- Mittels Silhouettenkoeffizient
 - hier wird die Nähe der einzelnen Datenpunkte zum eigenen Cluster im Gegensatz zu der Nähe der einzelnen Datenpunkte zum jeweiligen nächsten Nachbarcluster ermittelt
 - ein Koeffizient nahe 1 besagt, dass die Zuteilungen vermutlich gut sind
 - ein Koeffizient nahe -1 besagt das Gegenteil

Die Entscheidung liegt dann beim Betrachter!

kMeans – Das optimale k

```
1 from sklearn.datasets import load_iris
2 from sklearn.cluster import KMeans
3 from sklearn.metrics import silhouette_score
4 import matplotlib.pyplot as plt
5
6 iris = load_iris()
7 X = iris.data
8
9 # Das optimale k ermitteln
10 elbow, ss = [], []
11 kMin, kMax = 2, 10
12
13 for i in range(kMin, kMax+1):
14     kM = KMeans(n_clusters=i)
15     y_pred = kM.fit_predict(X)
16     silhouette_avg = silhouette_score(X, y_pred)
17     ss.append(silhouette_avg)
18     elbow.append(kM.inertia_)
19
20 # Die Ergebnisse ausgeben
21 fig = plt.figure(figsize=(14, 7))
22 fig.add_subplot(121)
23 plt.plot(*args: range(kMin, kMax+1), elbow, 'b-', label='Sum of squared distances')
24 plt.xlabel("Number of cluster")
25 plt.ylabel("SSD")
26 plt.legend()
27 fig.add_subplot(122)
28 plt.plot(*args: range(kMin, kMax+1), ss, 'b-', label='Silhouette Score')
29 plt.xlabel("Number of cluster")
30 plt.ylabel("Silhouette Score")
31 plt.legend()
32 plt.show()
```



Aufgrund der Ähnlichkeiten anhand der Merkmale wäre eine Einteilung in 2, 3 oder 4 Cluster denkbar.

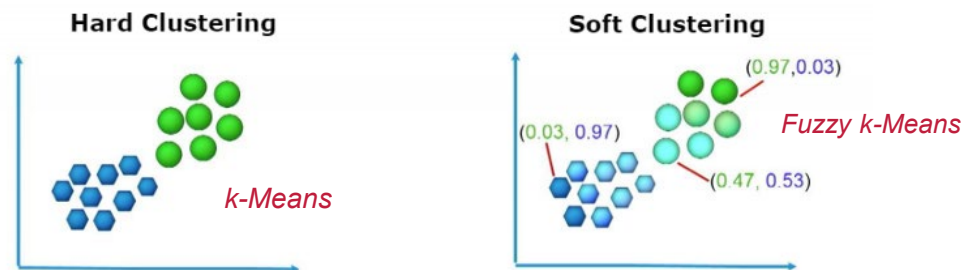
Das Ergebnis ist hier nicht eindeutig und weicht von der in diesem Beispiel bekannten Anzahl an Arten (3) ab.
(d.h. in diesem Beispiel könnten auch andere logische Einteilungen vorgenommen werden als die 3 bekannten Arten)

k-Means / k-Median – Fazit

- k-Median ist weniger anfällig für Ausreißer
- Vorteile:
 - Relativ einfache Implementierung
 - Einfache Anpassung an neue Daten (nach demselben Muster)
 - Skalierbar auf große Datensätze
- Nachteile:
 - das optimale k muss zuerst gefunden werden
 - Probleme bei Clustern unterschiedlicher Größe, Form und Dichte
 - „Fluch der Dimensionalität“

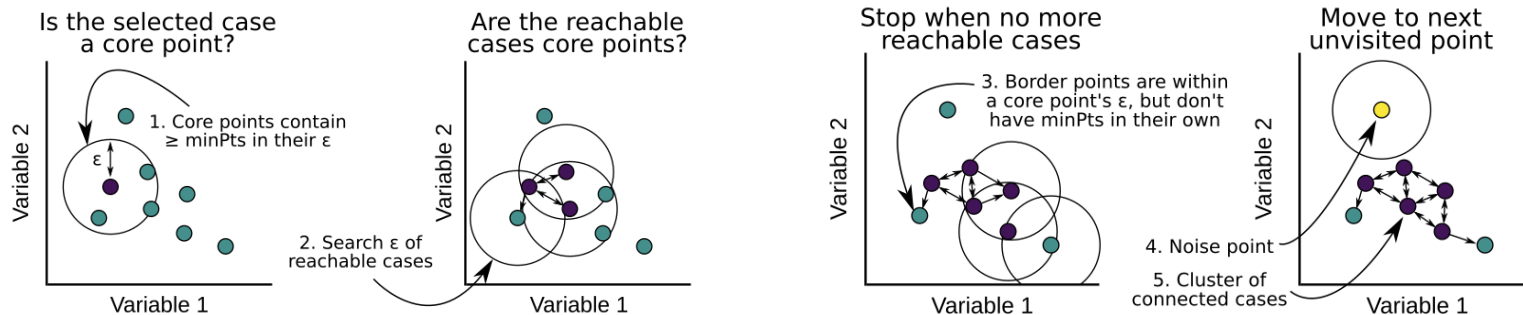
Fuzzy k-Means – Clusteranalyse

- Fuzzy Logik:
 - Unschärfelogik
 - Unterscheidung nach Gewichtung statt eindeutiger Zuordnung
- Erweiterung des k-Means – Algorithmus
- Die einzelnen Datenpunkte werden nicht genau einem Cluster zugeordnet, sondern mit Gewichten je Cluster versehen
- Diese Gewichte geben an, wie stark die Zugehörigkeit des Datenpunktes zum jeweiligen Cluster ist
- Entscheidung aufgrund der einzelnen Gewichte



DBScan – Clusteranalyse (unsupervised)

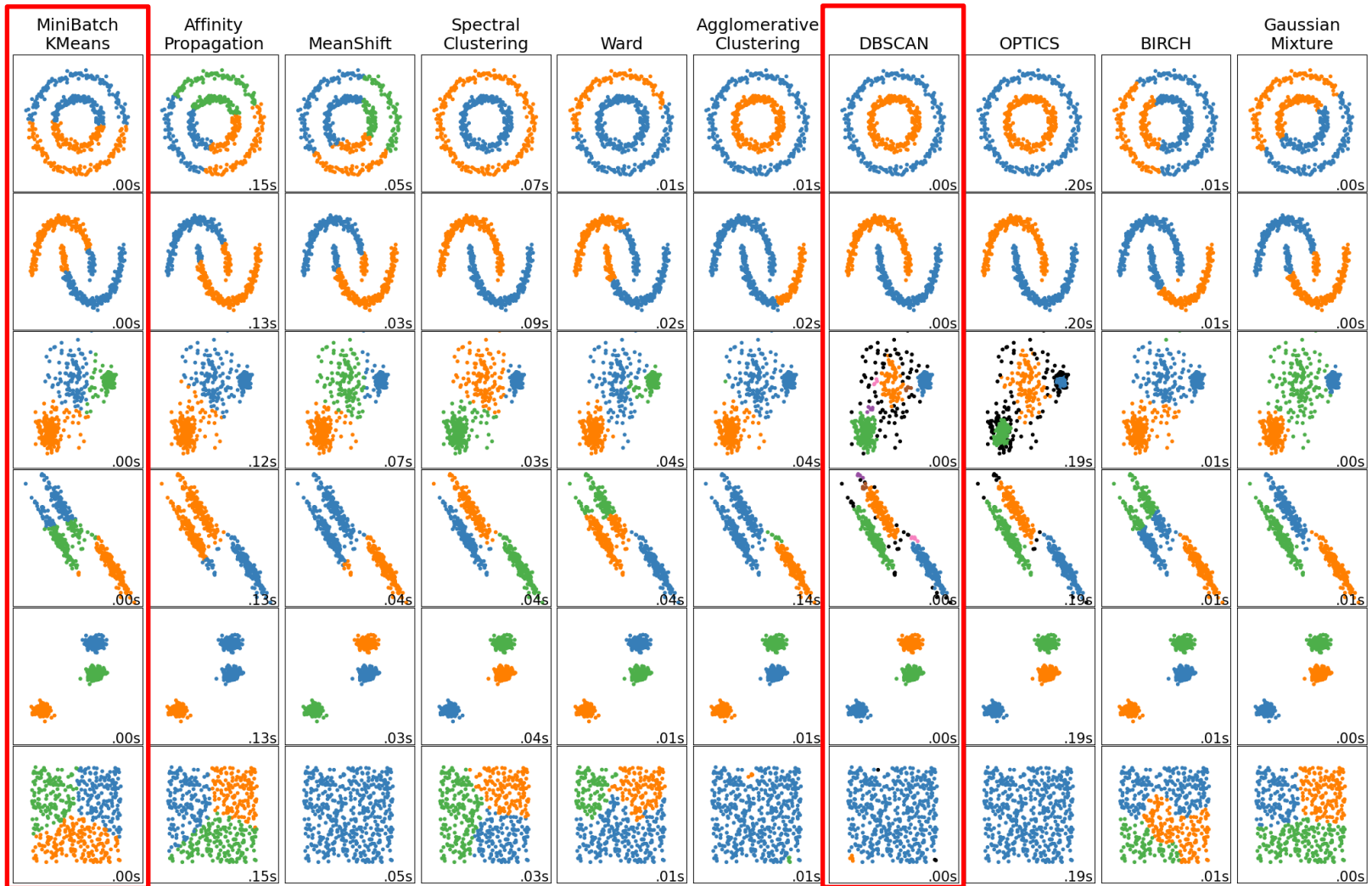
- Density-Based Spatial Clustering of Applications with Noise
- → “Dichtebasierte räumliche Clusteranalyse mit Rauschen“
- Zuordnung der einzelnen Datenpunkte nach Dichteverbundenheit
- Automatische Erkennung der Clusteranzahl
- Erkennung von Clustern beliebiger Form



- Parameter:

- Epsilon: max. Nachbarschaftslänge eines Datenpunktes
- minPts: min. Anzahl an Nachbarn

Vergleich der Cluster – Algorithmen (aus sklearn)



- [sklearn.neighbors.KNeighborsClassifier — scikit-learn 1.3.1 documentation](#)
- [sklearn.cluster.KMeans — scikit-learn 1.3.1 documentation](#)
- [Machine Learning 101 — Classification vs. Clustering | by Kevin C Lee | Medium](#)
- [K-Nearest neighbor clustering — Machine learning book \(vatsalparsaniya.github.io\)](#)
- [The 5 Clustering Algorithms Data Scientists Need to Know | by George Seif | Towards Data Science](#)
- [API Reference — scikit-learn 1.1.2 documentation](#)
- [Fuzzy C-Means Clustering —Is it Better than K-Means Clustering? | by Satyam Kumar | Towards Data Science](#)
- [DBSCAN Clustering in ML | Density based clustering – GeeksforGeeks](#)
- [Tutorial for DBSCAN Clustering in Python Sklearn - MLK - Machine Learning Knowledge](#)
- [18 Clustering based on density: DBSCAN and OPTICS - Machine Learning with R, the tidyverse, and mlr \(manning.com\)](#)
- [Understanding DBSCAN Algorithm and Implementation from Scratch | by Andrewngai | Towards Data Science](#)