

# HANSER



## Leseprobe

zum

## „Maschinelles Lernen“

von Jörg Frochte

ISBN (Buch): 978-3-446-45291-6

ISBN (E-Book): 978-3-446-45705-8

Weitere Informationen und Bestellungen unter  
<http://www.hanser-fachbuch.de/978-3-446-45291-6>

sowie im Buchhandel

© Carl Hanser Verlag, München

# Inhalt

<b>1</b>	<b>Einleitung</b> .....	<b>9</b>
<b>2</b>	<b>Maschinelles Lernen – Überblick und Abgrenzung</b> .....	<b>13</b>
	2.1 Lernen, was bedeutet das eigentlich? .....	13
	2.2 Künstliche Intelligenz, Data Mining und Knowledge Discovery in Databases ....	14
	2.3 Strukturierte und unstrukturierte Daten in Big und Small .....	17
	2.4 Überwachtes, unüberwachtes und bestärkendes Lernen .....	20
	2.5 Werkzeuge und Ressourcen .....	26
	2.6 Anforderungen und Datenschutz im praktischen Einsatz .....	27
<b>3</b>	<b>Python, NumPy, SciPy und Matplotlib – in a nutshell</b> .....	<b>32</b>
	3.1 Installation mittels Anaconda und die Spyder-IDE .....	32
	3.2 Python Grundlagen .....	35
	3.3 Matrizen und Arrays in NumPy .....	43
	3.4 Interpolation und Extrapolation von Funktionen mit SciPy .....	53
	3.5 Daten aus Textdateien laden und speichern .....	59
	3.6 Visualisieren mit der Matplotlib .....	61
	3.7 Performance-Probleme und Vektorisierung .....	65
<b>4</b>	<b>Statistische Grundlagen und Bayes-Klassifikator</b> .....	<b>68</b>
	4.1 Einige Grundbegriffe der Statistik .....	68
	4.2 Satz von Bayes und Skalenniveaus .....	70
	4.3 Bayes-Klassifikator, Verteilungen und Unabhängigkeit .....	76
<b>5</b>	<b>Lineare Modelle und Lazy Learning</b> .....	<b>88</b>
	5.1 Vektorräume, Metriken und Normen .....	88
	5.2 Methode der kleinsten Quadrate zur Regression .....	102
	5.3 Der Fluch der Dimensionalität .....	109
	5.4 k-Nearest-Neighbor-Algorithmus .....	110

<b>6</b>	<b>Entscheidungsbäume</b> .....	<b>117</b>
6.1	Bäume als Datenstruktur .....	117
6.2	Klassifikationsbäume für nominale Merkmale mit dem ID3-Algorithmus.....	122
6.3	Klassifikations- und Regressionsbäume für quantitative Merkmale .....	135
6.4	Overfitting und Pruning .....	149
6.5	Random Forest .....	154
<b>7</b>	<b>Ein- und mehrschichtige Feedforward-Netze</b> .....	<b>161</b>
7.1	Einlagiges Perzepton und Hebbsche Lernregel .....	162
7.2	Multilayer Perceptron und Gradientenverfahren .....	169
7.3	Auslegung, Lernsteuerung und Overfitting .....	189
<b>8</b>	<b>Deep Neural Networks mit Keras</b> .....	<b>210</b>
8.1	Deep Multilayer Perceptron und Regularisierung .....	210
8.2	Ein Einstieg in Convolutional Neural Networks .....	228
<b>9</b>	<b>Feature-Reduktion und -Auswahl</b> .....	<b>251</b>
9.1	Allgemeine Aufbereitung von Daten .....	253
9.2	Featureauswahl .....	261
9.3	Hauptkomponentenanalyse (PCA) .....	271
9.4	Autoencoder mit Keras .....	280
<b>10</b>	<b>Support Vector Machines</b> .....	<b>286</b>
10.1	Optimale Separation .....	286
10.2	Soft-Margin für nicht-linear separierbare Klassen.....	292
10.3	Kernel Ansätze .....	293
10.4	SVM in scikit-learn .....	298
<b>11</b>	<b>Clustering-Verfahren</b> .....	<b>304</b>
11.1	k-Means und k-Means++ .....	308
11.2	Fuzzy-C-Means .....	313
11.3	Dichte-basierte Cluster-Analyse mit DBSCAN .....	317
11.4	Hierarchische Clusteranalyse .....	324
<b>12</b>	<b>Bestärkendes Lernen</b> .....	<b>331</b>
12.1	Software-Agenten und ihre Umgebung .....	331
12.2	Markow-Entscheidungsproblem .....	334
12.3	Q-Learning .....	342

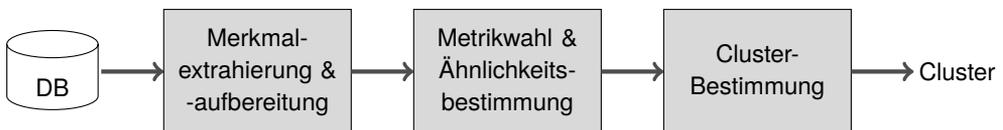
12.4 Der SARSA Algorithmus .....	349
12.5 Unvollständige Informationen und Softmax .....	351
12.6 Q-Learning mittels Funktionsapproximation .....	355
12.7 Ausblick auf Multi-Agenten- und hierarchische Szenarien .....	385
<b>Literatur .....</b>	<b>395</b>
<b>Index .....</b>	<b>401</b>

Clustering-Verfahren sind Verfahren für unüberwachtes Lernen. Das einzige andere Verfahren, das in diese Kategorie der unüberwachten Verfahren fällt, das wir besprochen haben, ist die PCA in Abschnitt 9.3, welche zur Merkmalsreduktion eingesetzt wird. Dass Clustering-Verfahren unüberwacht sind bedeutet, dass im Unterschied zur Klassifikation keine Zielwerte vorliegen. Während wir z. B. bei den Schwertlilien-Daten für alle 150 Datenbankeinträge wissen, um welche Art von Schwertlilien es sich handelt, haben wir dieses Wissen bei einem unüberwachten Verfahren nicht. Diese Spalte mit Daten würde hier fehlen und die Aufgabe unseres Algorithmus ist es dann, die Pflanzen danach zu gruppieren, welche Einträge sich am ähnlichsten sind.

Clusteranalyse ist somit die Einteilung einer Menge von Objekten in Gruppen, wobei wir folgende Ziele verfolgen:

- Wir wollen die Ähnlichkeit innerhalb der Gruppen maximieren.
- Wir wollen die Unterschiede zwischen den Gruppen maximieren, bzw. ihre Ähnlichkeit minimieren.

Weit präsenter als die Gruppierung von Pflanzen oder Tieren an Hand von Merkmalen ist den meisten Lesern sicherlich das Erlebnis im Bereich des Online-Handels. Wer ist nicht schon mal mit der Meldung konfrontiert worden: „Kunden, die diesen Artikel gekauft haben, kauften auch ...“. Es geht also oft um die Identifikation gleichartiger Käufergruppen- und Interessen. Vergleichbare Anwendungen gibt es auch in anderen Beziehungen wie Business-to-Business. Darüber hinaus haben die Verfahren ein Einsatzgebiet im Bereich des Information-Retrieval.

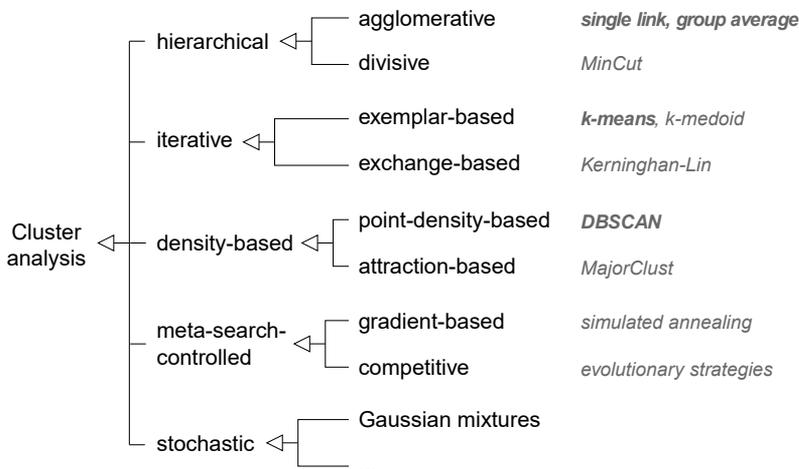


**Abbildung 11.1** Grundablauf einer Clusteranalyse

Der grundsätzliche Ablauf einer Clusteranalyse besteht aus drei Schritten, die in der Abbildung 11.1 skizziert sind. Das erste ist wie immer die Merkmalsextrahierung und -aufbereitung. Viele Dinge sind natürlich ähnlich wie schon in Kapitel 9 besprochen. Hierbei gibt es einen Unterschied zu den meisten überwachten Lernverfahren, die wir bis jetzt diskutiert haben. Bei der Analyse von Ähnlichkeiten ist es wie bei bereits besprochenen überwachten Verfahren wichtig, dass die verwendeten Einheiten in den Merkmalen das Verfahren nicht negativ beeinflussen. Nehmen wir einmal an, wir hätten bei dem bekannten Beispiel mit den Schwertlilien die Kelchblattlänge wie bisher in Zentimetern angegeben, die Kronblattlänge jedoch in Millimetern. Die meisten Algorithmen würden dann die Unterschiede bzgl. der Kronblätter völlig überbewerten und den Kelchblättern zu wenig Aufmerksamkeit schenken. Das ist im Wesentlichen das Problem, das wir bisher durch Skalierung oder Standardisierung gelöst haben und

auch in diesen Fall angehen müssten. Wichtiger für das Gelingen einer Clusteranalyse ist jedoch oft eher der Abstand der Elemente in einem Cluster zueinander. Nehmen wir an, es geht um die Flügelspanweite von Vögeln. Bei Wanderalbatrossen sind wir hier im Bereich von 350 Zentimetern und bei Tannenmeisen von eher 20 Zentimetern. Dazwischen liegt eine Zehnerpotenz, jedoch ist es wichtiger wie die Abstände innerhalb der Tannenmeisen aussehen, und dort wird oft ein Bereich von 17 bis 21 cm genannt. Der Spatz hingegen hat in der Regel eine etwas größere Spannweite um 23 cm. Wenn wir nun die Flügelspanweite auf das Intervall  $[0, 1]$  bringen würden ... was würde das für die Unterscheidbarkeit zwischen den kleinen Vogelarten bedeuten? Pauschal ist das schwer zu sagen und daher ist Domainwissen hier oft hilfreich, da man nur mit einer Vorerwartung sagen kann, ob man durch Standardisierung, Normierung etc. vielleicht in den Daten eher Unterschiede verwischt. In der Praxis liegt man allerdings mit einer Standardisierung der Daten oft richtig, jedoch sollte man die Möglichkeit hier Nachteile zu schaffen kritisch im Hinterkopf haben.

Es gibt natürlich sehr viele Clusterverfahren, die sich in unterschiedliche Einheiten gruppieren lassen. Die Abbildung 11.2 zeigt eine von Benno Stein vorgenommene Einteilungen von



**Abbildung 11.2** Taxonomie von Clusteralgorithmen

Cluster-Verfahren aus seinen Veranstaltungen [Ste18]. Wir werden aus den drei populärsten Gruppen, den hierarchischen, den iterativen und den dichte-basierten Verfahren, jeweils einen Clusteralgorithmus vorstellen. Dabei konzentrieren wir uns wieder jeweils auf recht populäre Verfahren wie den DBSCAN und k-Means.

Für unsere Cluster-Verfahren werden wir 4 Testfälle im Zweidimensionalen betrachten, da sie sich gut visualisieren lassen.

Die ersten beiden Fälle bestehen aus vier Bällen, die unterschiedlich viele Elemente beinhalten können. Das nächste ist eine Datenmenge, die aus drei Bällen besteht, zwischen denen es einen dünner besetzten Übergangsbereich gibt. Solche Datensätze werden oft als **Mouse-Dataset** bezeichnet. Nun folgen die schon aus Abschnitt 9.4 bekannten zwei Halbmonde und als letzte Datenmenge mit Struktur kommen zwei ineinanderliegende Kreise. Die sechste Menge von Daten hat gar keine Struktur und soll zeigen, was passiert, wenn die Verfahren auf Rauschen angewendet werden. Wichtig ist, sich später vor Augen zu halten, dass die Art, wie

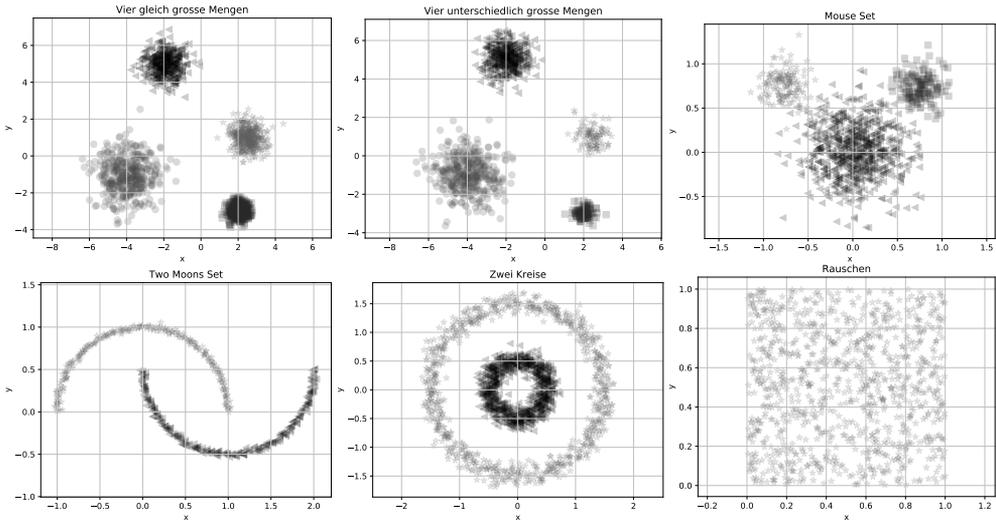


Abbildung 11.3 Testbeispiele für die Clusteralgorithmen

die Daten erzeugt wurden, nicht automatisch die einzige mögliche Art ist, diese als Cluster zu gruppieren. Als Menschen sehen wir Dinge anders als sie u. U. sind. Nehmen wir einmal Max Wertheimer, einen der Begründer der Gestaltpsychologie, und seine 1923 formulierten sechs wesentlichen Faktoren für Zusammenhänge in der Wahrnehmung als Grundlage. Menschen bevorzugen bei dem Beispiel mit den zwei Monden eben diese Zuordnung in zwei Halbkreise, da diese u. a. dem Gesetz der Nähe und der guten Gestalt entsprechen. Das ist jedoch nicht die einzige Gruppenbildung, die denkbar ist. Immerhin haben die beiden Enden der Bögen kaum etwas gemeinsam und liegen deutlich näher an Elementen des jeweils anderen Bogens. Ähnliches gilt für die beiden Kreise. Durch Farbe und Struktur käme kaum jemand auf die Idee, andere Gruppen als zwei Kreise zu bilden, aber vielleicht könnte man ja auch an der y-Achse die Mengen teilen und so zwei Gruppen erzeugen? Wir werden also gleich kritisch hinterfragen, warum uns ein Algorithmus welches Ergebnis zurückliefert. Manchmal ist es eine Schwäche der Methode, wenn das Ergebnis nicht unseren Erwartungen entspricht, jedoch ist unsere Erwartung auch manchmal zu stark von unseren Anlagen geprägt und rein formal nicht die einzige sinnvolle Lösung.

Die Quelltexte, um diese Beispiele zu erhalten, sind unten angegeben. Die Codezeilen zum Plotten der Daten sind nur für den ersten Fall angegeben, da diese sich in leichten Variationen immer wiederholen:

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def bubbleSetNormal(mx,my,number,s):
5     x = np.random.normal(0, s, number) + mx
6     y = np.random.normal(0, s, number) + my
7     return(x,y)
8
9 def fourBalls(n1,n2,n3,n4):
10    np.random.seed(42)
11    dataset = np.zeros( (n1+n2+n3+n4,2) )

```

```

12     (dataset[0:n1,0],dataset[0:n1,1]) = bubbleSetNormal( 2.5, 1.0,n1,0.5)
13     (dataset[n1:n1+n2,0],dataset[n1:n1+n2,1]) = bubbleSetNormal( 2.0,-3.0,n2,0.3)
14     (dataset[n1+n2:n1+n2+n3,0],dataset[n1+n2:n1+n2+n3,1]) = bubbleSetNormal(-2.0, 5.0,n3,0.6)
15     (dataset[n1+n2+n3:n1+n2+n3+n4,0],dataset[n1+n2+n3:n1+n2+n3+n4,1]) = bubbleSetNormal
    (-4.0,-1.0,n4,0.9)
16     return (dataset)
17
18 plt.close('all')
19 n1=n2=n3=n4=400
20 # Testbeispiel 1
21 XBalls = fourBalls(n1,n2,n3,n4)
22
23 fig = plt.figure(1)
24 ax = fig.add_subplot(1,1,1)
25 ax.scatter(XBalls[0:n1,0] ,XBalls[0:n1,1],c='red',s=60,alpha=0.2,marker='*')
26 ax.scatter(XBalls[n1:n1+n2,0],XBalls[n1:n1+n2,1],c='blue',s=60,alpha=0.2,marker='s')
27 ax.scatter(XBalls[n1+n2:n1+n2+n3,0],XBalls[n1+n2:n1+n2+n3,1],c='black',s=60,alpha=0.2,marker='
    <')
28 ax.scatter(XBalls[n1+n2+n3:n1+n2+n3+n4,0],XBalls[n1+n2+n3:n1+n2+n3+n4,1],c='green',s=60,alpha
    =0.2,marker='o')
29 ax.set_xlabel('x')
30 ax.set_ylabel('y')
31 ax.grid(True,linestyle='-',color='0.75')
32 ax.set_aspect('equal', 'datalim')
33 ax.set_title("Vier gleich grosse Mengen")
34
35 n1=n2=100
36 n3=n4=400
37 # Testbeispiel 2
38 XBalls2 = fourBalls(n1,n2,n3,n4)

```

Nun folgt der das *Mouse Dataset*. Es gibt da im Netz und in der Literatur unterschiedliche Varianten von. Teilweise gehen die Ohren noch stärker in den Kopf über. Wenn sie an dem Fall interessiert sind, können Sie das durch die Parameter in *bubbleSetNormal* leicht erreichen.

```

52 def mouseShape():
53     np.random.seed(42)
54     dataset = np.zeros( (1000,2) )
55     (dataset[0:150,0],dataset[0:150,1]) = bubbleSetNormal(-0.75, 0.75,150,0.15)
56     (dataset[150:300,0],dataset[150:300,1]) = bubbleSetNormal( 0.75, 0.75,150,0.15)
57     (dataset[300:1000,0],dataset[300:1000,1]) = bubbleSetNormal( 0, 0,700,0.29)
58     return (dataset)
59 # Testbeispiel 3
60 XBMouse = mouseShape()

```

Noch einmal der bekannte Fall mit den beiden Halbmoden, hier jedoch für ein Aufgabe beim unüberwachten Lernen.

```

73 def twoMoonsProblem( SamplesPerMoon=240, pNoise=2):
74     np.random.seed(42)
75     tMoon0 = np.linspace(0, np.pi, SamplesPerMoon)
76     tMoon1 = np.linspace(0, np.pi, SamplesPerMoon)
77     Moon0x = np.cos(tMoon0)
78     Moon0y = np.sin(tMoon0)
79     Moon1x = 1 - np.cos(tMoon1)
80     Moon1y = 0.5 - np.sin(tMoon1)
81     X = np.vstack((np.append(Moon0x, Moon1x), np.append(Moon0y, Moon1y))).T
82     X = X + pNoise/100*np.random.normal(size=X.shape)

```

```

83     Y = np.hstack([np.zeros(SamplesPerMoon), np.ones(SamplesPerMoon)])
84     return X, Y
85     # Testbeispiel 4
86     (XMoons,_) = twoMoonsProblem()

```

Nun zwei die zwei ineinander liegenden Kreise.

```

98     def circels():
99         np.random.seed(42)
100        phi = np.linspace(0,2*np.pi, 800)
101        x1 = 1.5*np.cos(phi)
102        y1 = 1.5*np.sin(phi)
103        x2 = 0.5*np.cos(phi)
104        y2 = 0.5*np.sin(phi)
105        X = np.vstack((np.append(x1,x2), np.append(y1,y2))).T
106        X = X + 0.1*np.random.normal(size=X.shape)
107        return(X)
108        # Testbeispiel 5
109        Xcircels = circels()

```

Abschließend noch das Rauschen im Einheitsquadrat.

```

121    np.random.seed(42)
122    # Testbeispiel 6
123    XRauschen = np.random.random( (1000,2) )

```

Beachten Sie, dass oben zwei unterschiedliche Arten zum Erzeugen von Zufallszahlen verwendet wurden. `np.random.random` erzeugt gleichmäßige Zufallszahlen zwischen 0 und 1, während `np.random.normal` entsprechend einer Normalverteilung diese erzeugt. Große Zufallszahlen sind also unwahrscheinlicher als keine, jedoch nicht ausgeschlossen. Daher enthalten die Mengen auch Einträge, die der eine oder andere Algorithmus als Ausreißer interpretieren wird.

## ■ 11.1 k-Means und k-Means++

Ziel von k-Means ist es, die Elemente  $x_j$  des Datenbestandes in  $k$  Cluster aufzuteilen. Das Besondere bei dieser Clustertechnik ist, dass es einen Repräsentanten  $\mu_i$  pro Cluster  $C_i$  gibt. Er ist quasi der Prototypvertreter des Clusters, jedoch im Regelfall kein Datenbankeintrag, sondern ein Mittelwert der Daten des Clusters. Die Frage ist nun, nach welchen Gesichtspunkten die Aufteilung erfolgt. Das Kriterium für die Qualität dieser Aufteilung ist, dass die Summe der Abweichungen von den Cluster-Repräsentanten in der gewählten Metrik minimal ist. Mathematisch entspricht dies der Optimierung der Funktion, hier bzgl. eines Minimums,

$$J = \sum_{i=1}^k \sum_{x_j \in S_i} d(x_j, \mu_i)$$

wobei  $d(x_j, \mu_i)$  die Distanz in der entsprechenden Metrik ist. Die häufigste Interpretation erfolgt jedoch über den **Lloyd-Algorithmus**, wenn wir statt einer beliebigen Metrik das Quadrat der Euklid-Norm als Grundlage für den Abstand spezialisieren, dann erhalten wir mit

$$d(x_j, \mu_i) = \|x_j - \mu_i\|^2$$

$$J = \sum_{i=1}^k \underbrace{\sum_{x_j \in C_i} \|x_j - \mu_i\|^2}_{(C)}. \quad (11.1)$$

Die äußere Summe durchläuft alle Cluster. Die Summe (C) besteht dabei aus dem Abstand aller Elemente des Clusters  $C_i$  vom Repräsentanten.  $J$  hängt hier nun mit der Methode der kleinsten Fehlerquadrate zusammen und man kann das Ziel im Sinne einer Varianzminimierung interpretieren. Unabhängig davon, welche Metrik man einsetzt, ist der Algorithmus in drei Phasen aufgeteilt.

1. Initialisiere  $k$  Repräsentanten  $\mu_i$  für die Cluster
2. Ordne jedes Element dem Cluster zu, bei welchem die Distanz zum Repräsentanten des Clusters am kleinsten ist
3. Berechne durch Mittelwertbildung die neuen Repräsentanten  $\mu_i$  der Cluster

Wenn der Schritt Nummer 3 abgeschlossen ist, wird wieder zu Nummer 2 gesprungen, bis die Lage der Repräsentanten stabil ist.

Der folgende Pseudocode des Algorithmus lässt sich leicht nach Python aber auch MATLAB oder R überführen.  $D$  ist dabei der gesamte Datenbestand, der vorliegt, und  $v$  ein Element daraus.

```

1: t = 0
2: for i = 1 to k do  $\mu_i = \text{choose}(D)$ 
3: end for
4: repeat
5:   t = t + 1
6:   for i = 1 to k do  $C_i = \emptyset$ 
7:   end for
8:   for all  $v \in D$  do
9:      $i = \text{argmin}_{j \in \{1..k\}} d(\mu_j, v)$ 
10:     $C_i = C_i \cup \{v\}$ 
11:   end for
12:   for i = 1 to k do
13:      $\mu_i^{\text{old}} = \mu_i$ 
14:      $\mu_i = \text{Mittelwert}(C_i)$ 
15:   end for
16:    $e = \sum_{i=1}^k d(\mu_i^{\text{old}}, \mu_i)$ 
17: until  $e < \epsilon$  OR  $t > t_{\text{max}}$ 

```

Die Laufzeit von k-Means ist  $O(n \cdot k \cdot r \cdot i)$ , wobei  $n$  die Anzahl der Einträge in unserer Datenbank  $D$  ist.  $r$  steht für die Dimension des Vektorraums, also die Anzahl der Merkmale.  $k$  ist die Zahl der Cluster, die k-Means finden soll, und  $i$  die Anzahl der Iterationen, welche bis zur Konvergenz benötigt werden. Für sinnvolle Einsatzfälle ist  $i$  dabei eine eher kleine Zahl im Bereich unter Hundert. Da für eine feste Fragestellung  $d$  konstant ist und ebenso die Anzahl der gesuchten Cluster, wird der Algorithmus oft als Methode mit einer linearen Komplexität bezeichnet, da er abhängig von der Datenbankgröße  $n$  im Wesentlichen nur linear wächst.

Im Allgemeinen arbeiten Clusteralgorithmen auf einer Gesamtmenge und teilen diese in einzelnen Cluster auf. Kommt ein neues Element in einer Datenbank hinzu, muss die Cluster-

analyse auch erneut durchgeführt werden. Bei k-Means gibt es jedoch die Repräsentanten, die man dazu nutzen kann, eine Vorhersage zu treffen, in welchem Cluster ein neues Element wohl wäre. Hierzu wird die Distanz zu jedem Repräsentanten berechnet und dann die Clusterzuordnung analog zur Zeile 9 in Pseudocode zurück geliefert. Daher können wir hier ausnahmsweise wie bei den überwachten Verfahren die Methoden `fit` und `predict` zur Verfügung stellen.

Wir beginnen mit der Methode, um ein Objekt zu instanziiieren. Wir beschränken uns bei unserer Umsetzung darauf,  $p$ -Normen für die Distanz zu verwenden, wobei der Default die Euklid-Norm mit  $p = 2$  ist. Notwendiger Parameter hingegen ist die Anzahl der gewünschten Cluster. Zusätzlich initialisieren wir noch eine Liste, um uns die Veränderung der Repräsentanten während des Algorithmus ansehen zu können.

```

1  import numpy as np
2
3  class kmeans:
4      def __init__(self, noOfClusters ,p=2):
5          self.noOfClusters = noOfClusters
6          self.p = p
7          self.fitHistory = []

```

Bei k-Means geht es darum Abstände von allen Zentren bzw. Repräsentanten der Cluster zu berechnen. Die Methode unten erlaubt dies in einer möglichst vektorisierten Form. Der erste Cluster wird dabei als Spezialfall vor der Schleife berechnet und abschließend alle weiteren Distanzen mittels `vstack` hinzugefügt.

```

8
9  def _computeDistances(self,X,centres):
10     distances = ( np.sum( np.abs(X-centres[0,:])**self.p,axis=1) )**(1/self.p)
11     for j in range(1,self.noOfClusters):
12         distancesAdd = ( np.sum( np.abs(X-centres[j,:])**self.p,axis=1) )**(1/self.p)
13         distances = np.vstack( (distances,distancesAdd) )
14     return(distances)

```

Die `fit`-Methode setzt mit dieser Distanzberechnung den oben angegebenen Pseudocode um.

```

15
16  def fit(self, X, maxIterations=42):
17      Xmin = X.min(axis=0)
18      Xmax = X.max(axis=0)
19      newcentres = np.random.rand(self.noOfClusters,X.shape[1])*(Xmax - Xmin)+Xmin
20      oldCentres = newcentres + 1
21      count = 0
22      while np.sum(np.abs(oldCentres-newcentres))!= 0 and count<maxIterations:
23          count = count + 1
24          oldCentres = np.copy(newcentres)
25          self.fitHistory.append(newcentres.copy())
26          distances =self._computeDistances(X,newcentres)
27          cluster = distances.argmax(axis=0)
28
29          for j in range(self.noOfClusters):
30              index = np.flatnonzero(cluster == j)
31              if index.shape[0]>0:
32                  newcentres[j,:] = np.sum(X[index,:],axis=0)/index.shape[0]
33              else:
34                  distances = ( np.sum( (X-newcentres[j,:])**self.p,axis=1) )**(1/self.p)
35                  i = distances.argmax(axis=0)

```

```

36         newcentres[j,:] = X[i,:]
37         cluster[i]=j
38     self.centres = newcentres
39     return(newcentres,cluster)

```

Auffällig hierbei ist der `else`-Zweig von Zeile 34 bis 37. Hierbei geht es um einen unangenehmen Spezialfall, in welchem ein Cluster leer bleibt. Tatsächlich ist es nämlich nicht ausgeschlossen, dass im Laufe einer Iteration einem Repräsentanten alle *Mitglieder* durch andere Repräsentanten abspenstig gemacht werden. Ist dies der Fall, wird der Eintrag in der Datenbank herausgesucht, der am nächsten an diesem Zentrum von Garnichts liegt, und dieser Eintrag wird dem leeren Cluster zugewiesen. Da er nur aus diesem einen Element besteht, wird dann auch der Repräsentant gleich dem einzigen Eintrag gesetzt.

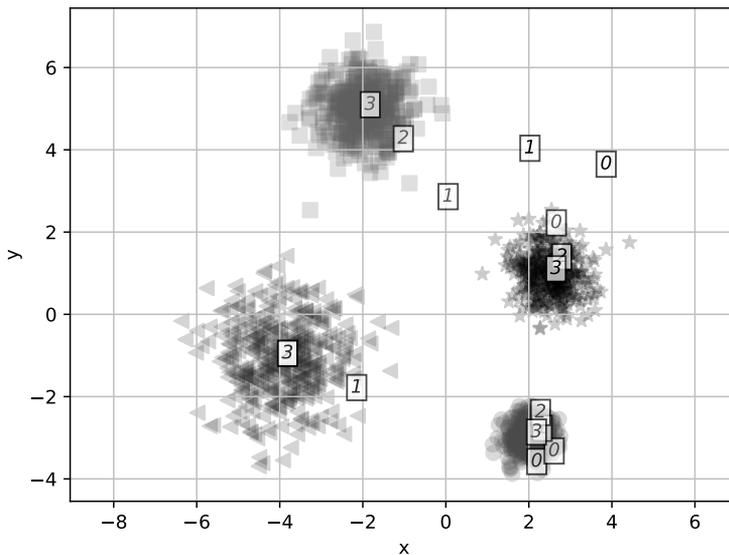
Wie schon erwähnt, können wir für k-Means eine `predict`-Methode bereitstellen:

```

40
41     def predict(self,X):
42         distances =self._computeDistances(X,self.centres)
43         cluster = distances.argmax(axis=0)
44         return cluster

```

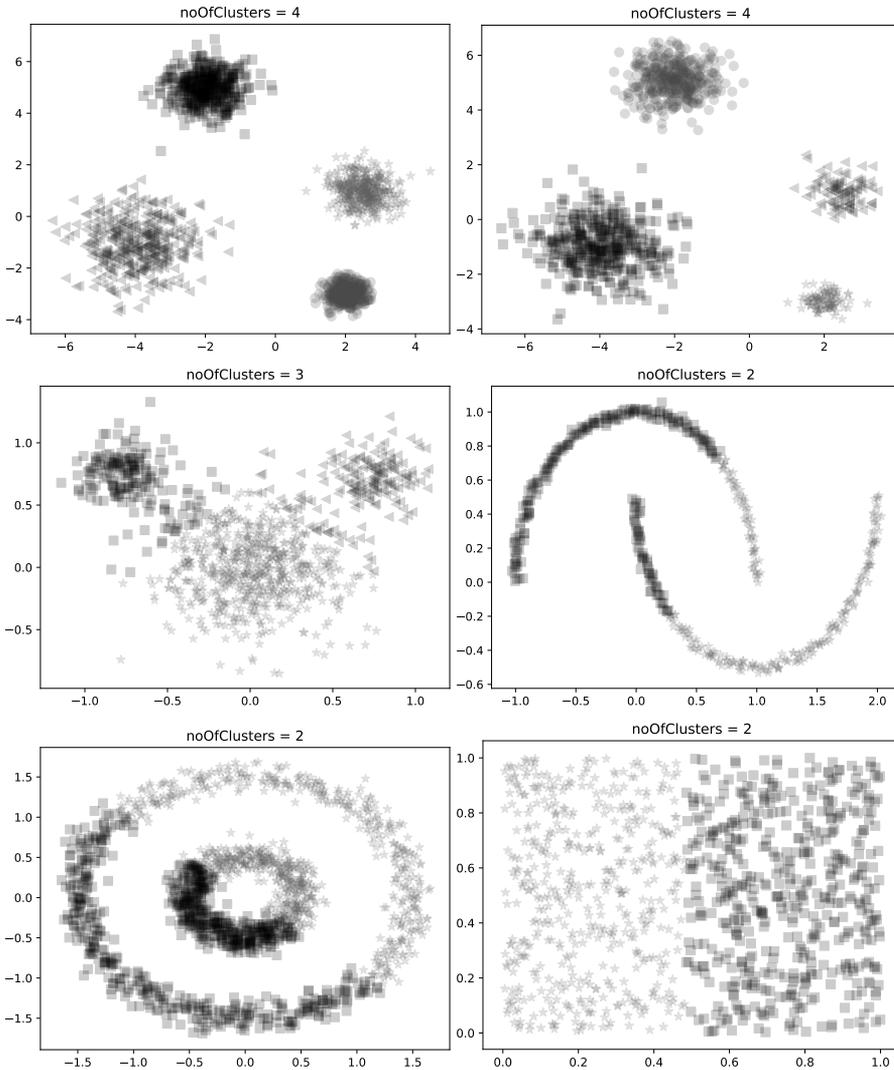
Nun testen wir unsere Implementierung auf dem Testfall 1, bei dem wir es mit vier Gruppen ähnlicher Mächtigkeit zu tun haben. Wie man in Abbildung 11.4 sehen kann, konvergiert das



**Abbildung 11.4** Testbeispiel 1 für k-Means

Verfahren sehr schnell. Nach nur vier Iterationen stehen die Repräsentanten und damit auch die einzelnen Cluster fest. Der Algorithmus war dabei auch in der Lage, mit etwas unglücklichen Startwerten umzugehen. Wenn man einmal sucht wo die vier mit 0 angegebenen Startpunkte liegen, so sind zwei nah beieinander in der grünen Menge gestartet. Auch die beiden später zum roten (Dreieck nach unten) und schwarzen (Kreis) Cluster gehörigen Startwerte sind eher eng beieinander. Während sich die nach oben und nach links zeigenden Dreiecke als Gruppen schnell getrennt haben, wurden die Kreise und nach unten zeigenden Dreiecke

zunächst in die gleiche Richtung gezogen und haben sich dann ab 2 deutlich aufgeteilt. Die unterschiedliche Dichte der Gruppen hatte hier keinen Einfluss auf das Clustering.



**Abbildung 11.5** Clusterbildung mit k-Means auf den Testproblemen

Die restlichen in Abbildung 11.5 dargestellten Ergebnisse auf den Testproblemen sind auch zufriedenstellend. Beim Mouse-Dataset sieht man, dass k-Means Gruppen gleicher Größe bevorzugt, weshalb die Ohren etwas auf den mittleren Bereich übergreifen und drei etwa gleich große Gruppen erzeugt werden. Die Testfälle 4 und 5 sehen für menschliche Betrachter etwas unglücklich aus, sind jedoch bzgl. der Minimierung des oben erwähnten Funktionals (11.1) stimmig. Eine für Menschen vertrautere Aufteilung erfolgt eher durch dichtebasierte Ansätze, wie wir sie in Abschnitt 11.3 besprechen werden. Wichtig ist jedoch vor allem der Testfall 2. Hier sind zwar die richtigen vier Cluster entstanden, jedoch ist dieses Ergebnis nicht stabil. Bei

unglücklichen Startwerten können auch die beiden rechten Gruppen zusammengelegt werden und eine der beiden größeren Gruppen stattdessen geteilt. Hier ist der Algorithmus also sehr abhängig von den Startwerten.

Um hier robustere Startbedingungen zu schaffen, wurde eine Variante von k-Means, nämlich **kmeans++**, entworfen. Der Unterschied zu k-Means bezieht sich dabei ausschließlich auf die Initialisierung der Repräsentanten. Mit der Menge  $D$  aller Datenbankeinträge  $x$  lässt sich das Vorgehen in wenigen Punkten zusammenfassen:

1. Wähle einen Repräsentanten  $\mu_1 \in D$  mit einer uniformen Wahrscheinlichkeit
2. Berechne für jeden Eintrag  $x$  den Abstand  $D(x)$  zum nächstgelegenen bereits gewählten Repräsentanten
3. Wähle zufällig einen neuen Datenpunkt als neuen Repräsentanten. Hierbei nutzt man jedoch keine uniforme Wahrscheinlichkeit, sondern gewichtet diese proportional zu  $D(x)^2$
4. Kehre zu Schritt 2 zurück, bis  $k$  Repräsentanten gewählt wurden
5. Führe nun den bekannten k-Means durch.

Wie man sieht sind hier die Startwerte automatisch Elemente der Datenbank und nicht rein zufällige Werte in dem Bereich, der durch die Datenbank plausibel erscheint. Durch den anschließenden Algorithmus werden jedoch wieder Mittelwerte gebildet, weshalb erneut die Repräsentanten in der Regel nicht mit Einträgen aus der Datenbank übereinstimmen werden.



Nutzen Sie k-Means, um das u. a. schon aus Abschnitt 3.5 bekannte *Iris flower data set* in Cluster einzuteilen; natürlich ohne die Zuordnung zu der jeweiligen Klasse zu verwenden. Wie Sie sehen werden, ist es hier nicht möglich, drei Cluster zu bilden, die genau die drei Gruppen von Lilien umfassen. Der Grund ist, dass diese Aufgabe als überwachtes Problem mit Zieldaten recht einfach ist, als unüberwachte Fragestellung jedoch nicht.

## ■ 11.2 Fuzzy-C-Means

Während k-Means zeitlich im Bereich der späten fünfziger und sechziger Jahre aufkam, wurde etwa ein Jahrzehnt später eine Fuzzy-Variante in [Dun73] durch J.C. Dunn vorgestellt und anschließend weiterentwickelt. Dazu muss man sich klar machen, dass der Begriff einer Fuzzymenge erst Mitte der Sechziger durch Lotfi Zadeh [Zad65] geprägt wurde. Die Grundidee ist etwas, was sofort intuitiv einleuchtet. Über die meisten realen Dinge sind Aussagen wie *X ist Y* selten unbestreitbar wahr, sondern eher graduell richtig. Beispielsweise könnten wir nun eine Reihe von Flussnamen aufzählen mit der Aussage *X ist ein langer Fluß*. Bei dem Nil und dem Amazonas ist die Zustimmung vermutlich gesichert. In dem Vorort, in dem ich wohne, fließt ein Bach mit Namen Kerbecke. Ohne eine Suchmaschine nicht auffindbar, wenn man hier nicht wohnt und daher sicherlich bzgl. der Aussage oben zu verneinen. Was machen wir denn nun mit dem Rhein? 1232 km sind gegen die 6852 km des Nil jetzt deutlich kürzer. Komplett verneinen würde jedoch nicht nur Rheinländer verärgern, sondern vermutlich auch damit alle Flüsse in Deutschland disqualifizieren. Die Ausweg ist eben die Fuzzymenge. Während man in der klassischen Logik diese Aussage nur mit Wahr (1) und Falsch (0) beantworten

kann, können wir hier einer Aussage einen Wahrheitswert zwischen 0 und 1 zuweisen. Mit dieser Fuzzy-Aussage zum Wahrheitsgehalt geht auch die Fuzzy-Zugehörigkeit mit einer Menge einher. Der Rhein wird also vielleicht nur mit einem Wahrheitswert von z. B. 0.6 zur Menge der langen Flüsse gehören.

Nehmen wir an, dass wir  $n$  Datensätze haben und  $C$  Cluster, die wir bilden wollen. Da wir für jeden der  $n$  Datensätze somit  $C$  Werte benötigen, für die Aussagen in welchem Maße der Datensatz zu einem Cluster zugehörig ist, ergibt sich eine Matrix  $W \in \mathbb{R}^{C \times n}$ . Diese Matrix enthält Werte von 0 bis 1. Jeder Eintrag  $w_{ij}$  gibt somit den Grad an, dem sich der Datensatz  $j$  dem Cluster  $i$  zugehörig fühlt. Ansonsten bleiben alle Basisideen von k-Means erhalten und es geht uns auch hierbei um eine Minimierung der Funktion

$$J = \sum_{i=1}^C \sum_{j=1}^n (w_{ij})^m d(x_j, \mu_i).$$

Neben der eher kosmetischen Änderung, dass  $k$  jetzt  $C$  heißt, geht es natürlich um den Faktor  $w_{ij}^m$  und seine Interpretation. Zum einen fließt die oben erwähnte Fuzzy-Zugehörigkeit so in das Funktional ein. Zum anderen geht es aber um die Rolle des Parameters  $m \geq 1$ . Dieser oft **fuzzifier** genannte Parameter verändert, wie scharf die Zugehörigkeit zu Clustern gewertet wird. Für  $m = 1$  geht die Zuordnung unverändert ein, je größer  $m$  wird, desto stärker werden jedoch Werte kleiner Eins reduziert. Ein großes  $m$  führt also zu kleineren Werten in der Matrix  $W$ , was einer gleichmäßigeren Zugehörigkeit zu den Clustern entspricht und deshalb zu unschärferen Clustern führt. Wählt man  $m = 1$ , erhält man nachdem der Algorithmus konvergiert ist, scharfe Mitgliedschaften wie schon beim k-Means. Werte größer als 3 sind unüblich und wenig erfolgversprechend. Sollte kein spezieller Grund durch Expertenwissen vorliegen, wird daher in der Regel die Mitte, also  $m = 2$ , als Ansatz gewählt.

Mit einem festen  $m$  ist der Algorithmus zum Auffinden des Minimums nun wieder sehr analog zu dem von k-Means. Wir erhalten nur eine weitere Nebenbedingung. Bei k-Means hatten wir nur die Bedingung: Die Cluster sind nicht-leer.

Im Fuzzy-Ansatz ist das noch leichter zu erfüllen. Es genügt zu fordern:

$$\sum_{j=1}^n w_{ij} > 0 \text{ für alle } i = 1 \dots C$$

Die neue Nebenbedingung ist, dass die Summe der Zugehörigkeiten zu den Clustern für jeden Datensatz 1 ist, das bedeutet

$$\sum_{i=1}^C w_{ij} = 1 \text{ für alle } j = 1 \dots n. \quad (11.2)$$

Man könnte sagen, die eine Bedingung bezieht sich auf die Zeilen unserer Matrix und die andere auf die Spalten.

Zur Lösung des Minimierungsproblems gehen wir nun wie folgt vor:

1. Initialisiere  $k$  Repräsentanten  $\mu_i$  für die Cluster
2. Berechne für jedes Element bzgl. jedes Clusters ein Maß der Zugehörigkeit mittels:

$$w_{ij} = \frac{1}{\sum_{k=1}^C \left( \frac{\|x_i - \mu_j\|}{\|x_i - \mu_k\|} \right)^{\frac{2}{m-1}}} = \frac{1}{\|x_i - \mu_j\|^{\frac{2}{m-1}} \cdot \sum_{k=1}^C \|x_i - \mu_k\|^{\frac{-2}{m-1}}} \quad (11.3)$$

3. Berechne durch gewichtete Mittelwertbildung die neuen Repräsentanten  $\mu_i$  der Cluster

$$\mu_i = \sum_{j=1}^n \frac{(w_{ij})^m}{\underbrace{\sum_{j=1}^n (w_{ij})^m}_{=\omega_j}} x_j = \frac{1}{\sum_{j=1}^n (w_{ij})^m} \sum_{j=1}^n (w_{ij})^m x_j \quad (11.4)$$

Die Gewichte in der Gleichung (11.3) findet man in der Literatur üblicherweise in der linken Formulierung. Die anschließende Umformung ist weniger intuitiv, eignet sich jedoch für eine vektorisierte Umsetzung in Python besser.

Jedoch ist schon die linke Formulierung der Gewichte vielleicht nicht direkt intuitiv zugänglich und man fragt sich beim Lesen: was passiert hier eigentlich? Man nimmt den Abstand zu einem Zentrum  $\mu_j$  und teilt diesen durch die Summe der Abstände zu allen Zentren, um sicherzustellen, dass die Summe aller dieser Gewichte eben entsprechend (11.2) 1 ergibt.

Schauen wir uns zum besseren Verständnis zwei Beispiele an, beide mit  $m = 2$  und drei Clusterzentren. Stellen wir uns ein Gewicht vor, bei dem der Abstand zu  $\mu_j$  sehr klein ist, also  $\varepsilon$ , und der Abstand zu den beiden anderen jeweils 1. Was erhalten wir dann als Wert?

$$\frac{1}{\left(\frac{\varepsilon}{\varepsilon}\right)^2 + \left(\frac{\varepsilon}{1}\right)^2 + \left(\frac{\varepsilon}{1}\right)^2} \approx 1$$

Das stimmt mit der Intuition überein. Wenn der Datensatz fast direkt neben dem Zentrum liegt, sollte das Gewicht auch fast 1 sein. Nun noch der Fall, in dem der Datensatz zu allen Zentren den gleichen Abstand, sagen wir  $1/2$  hat. Hier ergibt sich mit

$$\frac{1}{\left(\frac{1/2}{1/2}\right)^2 + \left(\frac{1/2}{1/2}\right)^2 + \left(\frac{1/2}{1/2}\right)^2} = \frac{1}{3}$$

wie gewünscht eine Gleichverteilung für alle Zentren.

Die Gleichung (11.4) folgt demselben Muster. Die Zentren ergeben sich als gewichtete Summe der Datensätze. Die Gewichte bestehen hier aus den Maßen für die Zugehörigkeit, also  $w_{ij}$  und der Summe aller dieser Werte für diesen konkreten Cluster. Durch die Konstruktion entsteht auch wieder der Effekt, dass die Summe aller  $\omega_j$  für einen Cluster 1 ergibt. Die rechte Formulierung von (11.4) ist wieder etwas schöner für die vektorisierte Umsetzung.

Mit dieser Umsetzung beginnen wir nun auch direkt, wobei die ersten 21 Zeilen identisch mit der k-Means Implementierung von Seite 310 sind. Lediglich der Klassenname muss durch `fuzzyCmeans` ersetzt und in Zeile 16 eine neue Option hinzugefügt werden.

```
16 def fit(self, X, maxIterations=42, vareps=10**-3):
```

Die wesentlichen Änderungen liegen in der `fit`-Methode. Hier starten wir mit Zeile 22 und ändern die Bedingung bzgl. des Abbruchs von 0 auf den neuen Parameter `vareps`. Eine entsprechende Anpassung lohnt sich auch ggf. für den k-Means, hier ist es jedoch wesentlich weil sonst zu viele Iterationen erfolgen.

```
22 while np.sum(np.abs(oldCentres-newcentres)) > vareps and count<maxIterations:
23     count = count + 1
24     oldCentres = np.copy(newcentres)
25     self.fitHistory.append(newcentres.copy())
```

Ab Zeile 26 setzen wir nun die Gleichungen (11.3) und (11.4) um. Hierzu berechnen wir wie zuvor erneut die Abstände und erlauben wieder nur p-Normen bzw. -Metriken. Darüber hinaus sind die Umsetzungen der Formeln auf den häufigsten Fall von  $m = 2$  spezialisiert. Wie man im Listing unten sieht, können alle Berechnungen vektorisiert erfolgen. Wegen  $m = 2$  haben wir es überall mit den Quadraten des Abstandes zu tun, weshalb wir den Output unserer Methode in Zeile 28 zunächst quadrieren. Anschließend berechnen wir in Zeile 29 die Summe aus (11.3). In Zeile 30 wird dann der Nenner von (11.3) gebildet und in Zeile 31 das Reziproke. Damit haben wir die Gewichte aus (11.3) fertig berechnet und können zu (11.4) übergehen. In Zeile bilden wir die benötigte Summe der Quadrate der Gewichte.

```

26
27     distances =self._computeDistances(X,newcentres)
28     d2 = distances**2
29     d2Sum = np.sum(1/d2,axis=0)
30     W = d2*d2Sum
31     W = 1/W
32     WSum = np.sum(W**2,axis=1)

```

Nun ist es unser Ziel, die  $\omega_j$ , die wir für jedes  $\mu_i$  berechnen, in einer einzigen  $i \times j$ -Matrix zu speichern, um so effektiver die Zentren zu berechnen. Dazu teilen wir die Matrix der Quadrate der Gewichte  $W$  durch die zuvor berechnete Summe. Damit die Dimensionen passen, muss transponiert werden. Anschließend kann jedoch die Verknüpfung von  $\omega_{ij}$  mit den Merkmalsvektoren direkt als Matrix-Vektor-Multiplikation durchgeführt werden, wodurch auch Schleifen zur Umsetzung von (11.4) vermieden werden können.

```

33     omega = ((W**2).T/WSum).T
34     newcentres = omega@X
35
36     self.fitHistory.append(newcentres.copy())
37     self.centers = newcentres
38     return(newcentres,W.T)

```

Am Schluss liefert die Methode die berechneten Zentren  $\mu_i$  in der Variablen `newcentres` zurück; darüber hinaus noch die Wahrscheinlichkeit, dass ein Element zu einem dieser Cluster gehört in Form der Matrix  $W$ . Eine `predict`-Methode auf der Basis dieser Zentren bzw. Repräsentanten kann wieder wie im Fall von k-Means umgesetzt werden.

Nimmt man nun den Fuzzy-C-Means und wendet ihn auf die Testprobleme an, so ergibt sich ein sehr ähnliches Bild wie für den k-Means auf Seite 11.5; jedenfalls wenn man einfach einen Datensatz dem Cluster zuschlägt, für den die größte Wahrscheinlichkeit vorliegt. Mittels

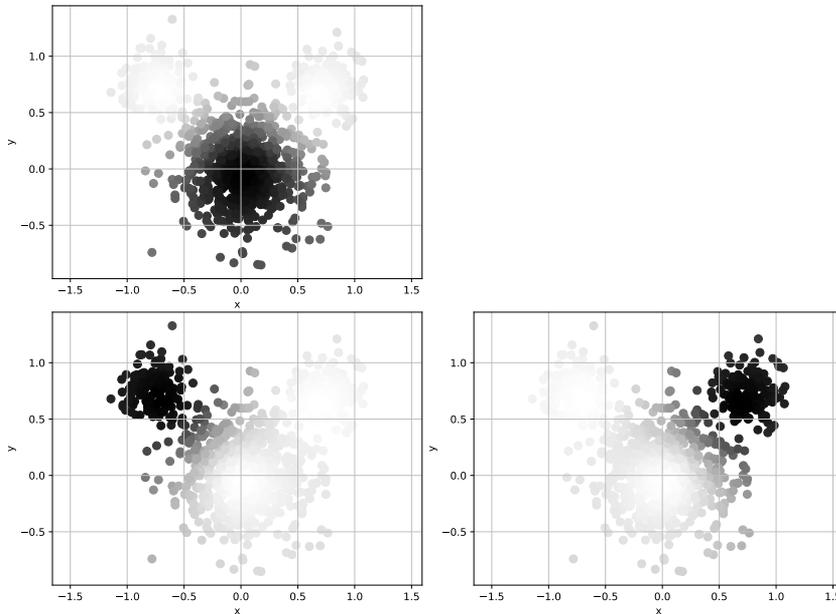
```

fig = plt.figure()
ax = fig.add_subplot(1,1,1)
ax.scatter(X[:,0] ,X[:,1],c=W[:,0],s=60, cmap='binary')

```

können wir uns den Grad der Zugehörigkeit zu einem Cluster visualisieren.

Wie man in Abbildung 11.6 sieht, ist dies nicht einfach eine Frage des Abstandes vom Repräsentanten. Die Ränder der *Ohren* zum Beispiel im ersten Bild sind dunkler. Der Grund ist, dass diese weniger sicher zu den anderen Clustern gehören und dadurch die Möglichkeit steigt, zum zentralen Cluster zu gehören. Der einzelne Punkt unten links hat wiederum zu keinem Cluster eine sehr hohe Zugehörigkeit. Diese Information des Fuzzy-Ansatzes kann man nutzen, um zum Beispiel Randbereiche abzutrennen oder die Gewichte in der späteren Verarbeitung von Daten neu zu bewerten. Werden die Gewichte hingegen nur im Sinn einer Zuordnung



**Abbildung 11.6** Wahrscheinlichkeit, zu einem speziellen Cluster zu gehören

über das Maximum verwendet, ist es in der Regel günstiger, den normalen k-Means zu verwenden. Die Ergebnisse sind in der Regel vergleichbar.



Wir haben im letzten Abschnitt bereits erlebt, dass die Clusteranalyse für das *Iris flower data set* nicht unproblematisch ist. Wie verhält sich die Fuzzy-Variante hier? Visualisieren Sie sich die Gebiete, in denen das Verfahren besonders sichere oder unsichere Zuordnungen macht, mit Hilfe mehrerer 2D- oder 3D-Plots.

## ■ 11.3 Dichte-basierte Cluster-Analyse mit DBSCAN

**DBSCAN** (Density-Based Spatial Clustering of Applications with Noise) gehört zur Klasse der dichte-basierten Clusteralgorithmen. Er wurde 1996 in der Veröffentlichung [EEHJ96] vorgestellt. Später wurde die Grundidee in Formen wie **OPTICS** (Ordering Points To Identify the Clustering Structure) [ABKS99] oder **MajorClust** [SN99] abgewandelt, wobei wir hier nur die Grundform von DBSCAN besprechen. Tatsächlich ist auch diese Grundform diejenige, die man in Bibliotheken und Veröffentlichungen am meisten verwendet findet. Dieser dichte-basierte Algorithmus kann mehrere Cluster erkennen, ohne dass wie bei den k-Means-Varianten zuvor die Anzahl der Cluster bekannt sein muss. Darüber hinaus werden Rauschpunkte im Laufe der Clusteranalyse erkannt, für das Clustering ignoriert und separat zurückgeliefert. Zu den Nach-

# Index

- $\epsilon$ -greedy, 344
- $\epsilon$ -greedy
  - Exploitation, 344
  - Exploration, 344
- Überanpassung, → Overfitting
- Überwachtes Lernen, 20
  
- Accuracy, 236
- Action-Value-Funktion, → Q-Funktion
- activation function, → Aktivierungsfunktion
- Agent
  - Actuators, 332
  - Condition-action rules, 332
  - Critic Modul, 334
  - goal-based agents, 331
  - kognitiv, 331
  - learning agents, 331
  - Learning Element, 333
  - model-based reflex agent, 331
  - Performance Element, 333
  - Problem Generator, 334
  - robust, 331
  - Sensors, 332
  - simple reflex agents, 331
  - sozial, 331
  - State, 332
  - utility-based agents, 331
- Agglomerative Clusterverfahren, 324
- Aktivierungsfunktion, 162
  - linear, 173
  - Rectifier Linear Unit, 211
  - Sigmoidfunktion, 171
  - Tangens Hyperbolicus, 171
- Anonymisierung, 29
- Average Pooling, 235
- Average-Linkage, 326
- Axiome von Kolmogorow, 69
  
- Backpropagation, 178
- Bagging, 154
  
- Basis, 94
  - kanonisch, 95
- Basiswechsel, 95
- Batch-Learning, 181
- Baumhöhe, 124
- bedingte Wahrscheinlichkeiten, 70
- Bellmansches Optimalitätsprinzip, 342
- Bestärkendes Lernen, 23, 331
  - Batch Reinforcement Learning, 369
  - Belohnung, 337
  - Environment, 332
  - Growing Batch Reinforcement Learnings, 369
  - kumuliertiver Reward, 339
  - modellbasiert, 342
  - modellfrei, 342
  - off-policy-Learning, 350
  - on-policy-Learning, 350
  - optimale Aktion, 341
  - optimale Policy, 340
  - Policy, 337
  - Reward, 337
  - Strategie, 337
  - Value Function, 339
- Bias, 82
- Bias-Neuron, 165
- Big Data, 18
- Boltzmann-Funktion, 353
- Boosting, 154
- Bootstrap-Sample, 155
  
- Centroid-Method, 326
- Classification Error, 236
- Cleverbot, 15
- CNN, 228
  - Bias-Neuron, 234
- Complete-Linkage, 325
- ConvNet, 228
- Convolutional Neural Networks, 228
- CPython, 65

- Cross-Correlation, 230
- Cross-Entropy, → Kreuzentropie
- Cross-Entropy-Error, 237
- Curse of Dimensionality, 19, 109
- CVXOPT, 286
  
- Data Augmentation, 249
- Dataset
  - Acute Inflammations Data Set, 73
  - Bevölkerungsentwicklung der BRD, 196
  - Bike Sharing Data Set, 145
  - Boston Housing Dataset, 102, 206, 212
  - CIFAR-10, 240
  - Fisher’s Iris Data Set, 59
  - MNIST, 63, 225
  - Mouse-Dataset, 305
  - Two Moons, 112
- Datensatz, 17
- Datenschutz, 29
- Datensicherheit, 29
- DBSCAN, 317
- Deep Autoencoder, 281
- Deep Networks, 174
- Dendrogramm, 324
- Dichteverbundenheit, 318
- Dimension des Vektorraums, 95
- discount factor, 338
- disjunkt, 71
- Diskontierungsfaktor, → discount factor
- Divisive Clusterverfahren, 324
- Dropout, 223
- Duale Formulierung, 291
- Duales Problem, 291
- Duck-Typing, 38
  
- Eager Learner, 23
- Eager Learning, 110
- early stopping, 199, 204
- einlagiges Perzeptron, 165
- Ensemble Learning, 154
- Entropie, 127
- Entscheidungsbaum
  - Aufteilung, 123
- equivariant representation, 232, 233
- Ereignis, 69
- Ereignisraum, 69
- Ergebnisraum, 68
  
- Ergebnisse, 69
- Erwartungswert, 81
- Erzeugendensystem, 94
- Experience Replay, 370
  - Mini-Batch, 371
- External Path Length, 124
- Extrapolation, 59, 196
  
- Feature, 17
- Feature Map, 229
- feedforward-Netze, 174
- Fehler
  - statistisch, 82
  - systematische, 82
- Flattening, 235
- Fluch der Dimensionalität, → Curse of Dimensionality
- Formel von Lance und Williams, 327
- Freiheitsgrade ein neuronales Netz, 189
- Fully-connected Layer, 212, 235
- Fully-connected Neural Network, 212
- fuzzifizier, 314
  
- Gauß-Newton-Verfahren, 102
- Gaußschen Glockenkurve, 81
- Generalisierung, 196
- geometrische Reihe, 338
- Gewichtete Pfadlängensumme, 124
- Gini Impurity, 136
- GMRES, 108
- GNU Octave, 46
- Gradient, 176
- Gramsche Matrix, 296
- Greedy-Algorithmus, 266
- Gridworld, 350
  
- Hauptkomponentenanalyse, 271
- Hebbsche Lernregel, 166
- Hiddenlayer, 172
- Hierarchie von Lernebenen, 389
- Histogramm, 83
- Hitchbot, 15
- Hyperebene, 103
  
- ImageNet, 247
- Imputation, 257
- Incremental Learning, 181
- Information Gain, 129
- Informationsgehalt, 126

- Input, 229
- Input-Layer, 172
- Intervallskala, 74
- IPython, 34
- Joint Action Learning, 386
- Jupyter-Notebook, 33
- k-d-Baum, 113
- k-Nearest-Neighbor-Algorithmus, 110
- k-NN, 110
- k-NN classification, 111
- k-NN regression, 111
- Künstliche Intelligenz, 14
  - schwache, 15
  - starke, 14
- Kaggle, 27
- Kardinalskala, 74
- Karush-Kuhn-Tucker-Bedingungen, 289
- Kausalität, 77
- Keras, 11, 210
  - AveragePooling2D, 242
  - Callbacks, 215
  - Conv2D, 242
  - early stopping, 214
  - EarlyStopping, 215
  - evaluate, 227
  - ImageDataGenerator, 249
  - Keras function, 245
  - ModelCheckpoint, 215
  - Sequential model, 212, 213
  - Trainable, 248
  - utils.to\_categorical, 226
- Kernel, 229, 296
- Kernel Matrix, 296
- Kernel Methods, 286, 296
- Klassifikation, 21
- Klassifizierungsproblem, 22
- kmeans++, 313
- Knowledge Discovery in Databases, 15
- Konfusionsmatrix, 80
- Koordinatenform, 103
- Korrelation, 77
- Korrelationskoeffizient, →
  - Pearson-Korrelationskoeffizient
- Korrelationsmatrix, 263
- Kovarianz, 261, 262
- Kovarianzmatrix, 273
- Kreuzentropie, 237
- Kreuzvalidierung, 86
  - Holdout, 86
  - k-fach, 86
  - Monte Carlo-Wiederholungen, 86
- L1-Regularisierung, 217
- L2-Regularisierung, 217
- Lagrange-Multiplikatoren, 289
- LAPACK, 108
- Lazy Learner, 23
- Lazy Learning, 110
- Least Squares Method, 106
- Lernhierarchie, 389
  - Aufgabe, 387
  - Plan, 387
  - Reaktion, 388
  - Reflexe, 389
  - Task, 387
- Lernrate, 182
- linear unabhängig, 94
- Lineare Separierbarkeit, 164
- linearer Kernel, 296
- Linearkombination, 92
- Lloyd-Algorithmus, 308
- Mahalanobis-Distanz, 99
- MajorClust, 317
- Manhattan-Metrik, 95
- Manhattan-Norm, 95
- Margin, 289
- Markov Decision Process, 337
  - deterministisch, 337
- Maslows Hammer, 10
- MATLAB, 46
- Matplotlib, 54
  - add\_subplot, 61
  - array, 54
  - Axes3D, 63
  - axis-Objekt, 61
  - cmap, 63
  - hold, 56
  - imshow, 64
  - pcolormesh, 66
  - plt, 56
  - pyplot, 56

- scatter, 61
- set\_xlabel, 61
- set\_ylabel, 61
- show, 56
- tight\_layout, 61
- Max-Pooling, 235
- Mean Square Error, 237
- Median, 254
- Mehrlagiges Perzeptron, 174
- Merkmal, 17
- Merkmalsausprägung, 78
- Metrik, 98
- Mittelwert der Stichprobe, 82
- Multilayer perceptron, → Mehrlagiges Perzeptron
  
- Nash Q-learning, 386
- natürlichen Nullpunkt, 74
- Neural Fitted Q Iteration, 372
- Neuronales Netz
  - fully connected, 174
  - Initialisierung der Gewichte, 183
  - Sättigung, 183
  - vollvermascht, 174
- Nominalskala, 73
- Norm, 96
- Normierung, 254
- NumPy, 44
  - arange, 46
  - array, 44
  - Array Broadcasting, 51
  - Array Slicing, 47
  - choice, 50
  - Deep Copy, 48
  - delete, 50
  - Erzeugen von Arrays, 46
  - Erzeugen von Vektoren, 46
  - flatnonzero, 49
  - ix\_, 49
  - loadtext, 60
  - maximum, 211
  - meshgrid, 65
  - poly1d, 55
  - rand, 50
  - randint, 50
  - ravel, 67
  - reshape, 47
  - seed, 51
  - set\_printoptions, 51
  - size, 49
  - View, 48
  - Zufallszahlen, 50
- Ockhams Rasiermesser, 55
- Offline-Learning, 182
- Offset, 103
- on-Neuron, 165
- Online-Learning, 182
- OpenAI Gym, 373
- OpenCV, 26
- OPTICS, 317
- Ordinalskalar, 73
- Orientierung, 357
- Out-of-Bag-Error, 156
- Outlier Detection, 26
- Outputlayer, 172
- Overfitting, 76, 200
  
- p-Normen, 96
- Pandas, 26
- parameter sharing, 232, 233
- Paris-Metrik, 99
- partially observable Markov decision process, 351
- Partielle Ableitung, 176
- Pearson-Korrelationskoeffizient, 78, 262
- Pfadlängensumme, 124
- Polygonzug, 54
- POMDP, → partially observable Markov decision process
- Pooling, 235
- Post-Pruning, 151, 152
- Pre-Pruning, 151
- Principal Component Analysis, 271
- Pruning, 151, 152
- Pseudonymisierung, 29
- Python
  - Built-in Types, 36
  - def, 38
  - Dictionaries, 36
  - doc String, 39
  - dynamische Typisierung, 38
  - import, 39
  - importlib, 40

- Introspection, 34
  - Lists, 36
  - None, 38
  - pretty printer, 51
  - Private Methode, 42
  - PYTHONPATH, 40
  - Referenzen, 36
  - return, 38
  - str, 39
  - Strings, 36
  - Tuples, 36
  - type, 35
- Q-Function, 342
- radial basis function, → Radiale Basisfunktionen
- Radiale Basisfunktionen, 296
- Random Forest, 155
- Rationalskala, 74
- Record, 17
- Rectifier Linear Unit, 211
- Recurrent Neural Network, 375
- Reduced-Error Ansatz, 152
- Regressionsproblem, 22
- Reinforcement Learning, → Bestärkendes Lernen
- ReLU
- see Rectifier Linear Unit, 211
- Residuenquadratsumme, 105
- RoboCup Simulation League, 27
- Runges Phänomen, 57
- Satz von Bayes, 70
- Satz von der totalen Wahrscheinlichkeit, 71
- Scatter Plots, 61
- Schlupfvariablen, 292
- scikit-learn, 11, 151
- DecisionTreeClassifier, 149
  - DecisionTreeRegressor, 149
  - GaussianNB, 87
  - Imputer, 259
  - KNeighborsClassifier, 116
  - KNeighborsRegressor, 116
  - LinearRegression, 109
  - LinearSVC, 298
  - MLPRegressor, 208
  - PCA, 280
  - RandomForestClassifier, 160
  - RandomForestRegressor, 160
  - RFE, 271
  - SVC, 298
- SciPy, 44
- dendrogram, 327
  - interpolate, 56
  - lagrange, 56
  - linkage, 327
  - pdist, 327
- Seaborn, 26
- Semi-überwachtes Lernen, 249
- Sequential Backwards Selection, 266
- Sequential Forward Selection, 270
- sequenzielle Rückwärtsauswahl, 266
- sequenzielle Vorwärtsauswahl, 270
- Sigmoidfunktion, 171
- Simulation Data Mining, 19
- Single-Linkage, 325
- Skalenniveaus, 73
- slack variables, → Schlupfvariablen
- Softmax, 237
- Softmax-Funktion, 353
- Softwarepatenten, 223
- sparse interactions, 232
- Spyder, 33
- %matplotlib, 64
  - Command History, 34
  - Introspection, 34
  - Tab Completion, 34
  - Variable Explorer, 33
- Stützstellen, 55
- Standardabweichung, 81
- empirische, 82
- Standardisierung, 255
- Stichprobenvarianz, 82
- Stochastic gradient descent, 182
- strukturierte Daten, 17
- Studentisierung, 255
- Subagging, 155
- Summe der Fehlerquadrate, 57, 105
- Support Vector Machines, 286
- one-vs-all, 297
  - one-vs-one, 297
  - one-vs-rest, 298
- Support Vectors, 289

- SVM, → Support Vector Machines
- SymPy, 26
- Tangens Hyperbolicus, 171
- TensorFlow, 210
- Testmenge, 75, 199
- Theano, 210
- time serie, → Zeitreihen
- Toeplitz Matrizen, 232
- TORCS, 27
- Trainingsmenge, 75, 199
- Transparenz, 28
- Turing-Test, 15
- UCI Machine Learning Repository, 27
- Unüberwachtes Lernen, 24
- Unity3D, 27
- unstrukturierte Daten, 18
- Unterraum
  - affinen, 286
- Validierungsmenge, 75, 151, 199
- Vektor von Bias-Neuronen, 213
- Verhältnisskala, 74
- Verzerrung, 82
- Wahrscheinlichkeit
  - gemessene, 69
- Weighted external path length, 124
- XGboost, 27
- YOLO, 249
- Zeitreihen, 228
- Zero Padding, 231
- Zufallsbeobachtung, 68
- Zufallsexperiment, 68
- Zurücklegen, 154