

---

# Entscheidungsbäume

*Entscheidungsbäume* sind sehr flexible Machine-Learning-Algorithmen, die sich sowohl für Klassifikations- als auch für Regressionsaufgaben eignen. Sogar Aufgaben mit multiplen Ausgaben lassen sich mit ihnen lösen. Es sind sehr mächtige Algorithmen, mit denen sich komplexe Datensätze fitten lassen. Beispielsweise haben Sie in Kapitel 2 ein Modell mit dem `DecisionTreeRegressor` trainiert und an den Datensatz zu kalifornischen Immobilien perfekt angepasst (genauer gesagt, overfittet).

Entscheidungsbäume sind außerdem die funktionelle Komponente von Random Forests (siehe Kapitel 7), die zu den mächtigsten heute verfügbaren Machine-Learning-Algorithmen gehören.

In diesem Kapitel werden wir besprechen, wie sich Entscheidungsbäume trainieren, visualisieren und für Vorhersagen einsetzen lassen. Anschließend werden wir den in Scikit-Learn verwendeten CART-Trainingsalgorithmus durchgehen. Wir werden betrachten, wie sich Entscheidungsbäume regularisieren und für Regressionsaufgaben einsetzen lassen. Schließlich werden Sie einige Einschränkungen von Entscheidungsbäumen kennenlernen.

## Trainieren und Visualisieren eines Entscheidungsbaums

Um Entscheidungsbäume zu verstehen, wollen wir zunächst einen erstellen und uns ansehen, wie dieser Vorhersagen trifft. Der folgende Code trainiert einen `DecisionTreeClassifier` auf dem Iris-Datensatz (siehe Kapitel 4):

```
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier

iris = load_iris(as_frame=True)
X_iris = iris.data[["petal length (cm)", "petal width (cm)"]].values
y_iris = iris.target

tree_clf = DecisionTreeClassifier(max_depth=2, random_state=42)
tree_clf.fit(X_iris, y_iris)
```

Sie können den trainierten Entscheidungsbaum visualisieren, indem Sie durch Aufrufen der Funktion `export_graphviz()` eine Datei namens `iris_tree.dot` mit einer Repräsentation als Graph erzeugen:

```

from sklearn.tree import export_graphviz

export_graphviz(
    tree_clf,
    out_file="iris_tree.dot",
    feature_names=["petal length (cm)", "petal width (cm)"],
    class_names=iris.target_names,
    rounded=True,
    filled=True
)

```

Dann können Sie `graphviz.Source.from_file()` nutzen, um die Datei in einem Jupyter Notebook zu laden und anzuzeigen:

```

from graphviz import Source

Source.from_file("iris_tree.dot")

```

Graphviz (<https://graphviz.org>) ist ein Open-Source-Softwarepaket zur Visualisierung von Graphen. Es enthält auch ein Befehlszeilentool `dot`, mit dem sich `.dot`-Dateien in eine Vielzahl von Formaten umwandeln lassen, zum Beispiel nach PDF oder PNG.

Ihr erster Entscheidungsbaum ist der in Abbildung 6-1 dargestellte.

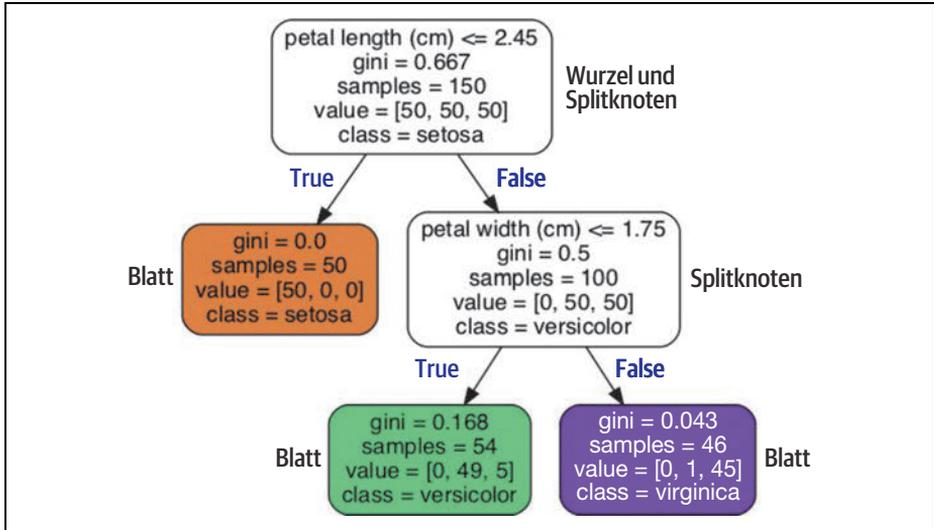


Abbildung 6-1: Iris-Entscheidungsbaum

## Vorhersagen treffen

Betrachten wir nun, wie der in Abbildung 6-1 dargestellte Baum Vorhersagen trifft. Nehmen wir an, Sie finden eine Iris-Blüte und möchten diese anhand der Kronblätter klassifizieren. Sie beginnen an der *Wurzel* des Baums (bei depth 0 ganz oben): Dieser Knoten stellt die Frage, ob das Kronblatt der Blüte kürzer als 2,45 cm ist. Ist das der Fall, fahren Sie mit dem Kind auf der linken Seite fort (depth 1, links). In diesem Fall ist der Knoten ein *Blatt* (es besitzt keine weiteren Kinder) und stellt keine weiteren Fragen: Sie übernehmen einfach die von diesem Knoten vorhergesagte Kategorie, somit sagt unser Entscheidungsbaum vorher, dass Ihre Blüte eine *Iris setosa* (class=setosa) ist.

Angenommen, Sie finden jetzt eine weitere Blüte, bei der das Kronblatt diesmal länger als 2,45 cm ist. Sie beginnen wieder an der Wurzel, fahren nun aber mit dem rechten Kind fort (bei depth 1, rechts). Dieser Knoten ist kein Blatt, sondern ein Splitknoten, der eine weitere Frage stellt: Ist das Kronblatt schmaler als 1,75 cm? Ist dies der Fall, ist Ihre Blüte vermutlich eine *Iris versicolor* (bei depth 2, links). Wenn nicht, ist sie vermutlich eine *Iris virginica* (bei depth 2, rechts). Es ist tatsächlich so einfach.



Einer der vielen Vorzüge von Entscheidungsbäumen ist, dass sie sehr wenig Vorbereitung der Daten erfordern. Insbesondere ist keinerlei Skalierung oder Zentrierung von Merkmalen notwendig.

Das Attribut *samples* eines Knotens zählt, für wie viele Trainingsdatenpunkte dieser gültig ist. Beispielsweise haben 100 Datenpunkte ein Kronblatt mit einer Länge von mindestens 2,45 cm (bei depth 1, rechts), davon haben 54 ein Kronblatt mit einer Breite von weniger als 1,75 cm (bei depth 2, links). Das Attribut *value* eines Knotens verrät uns, wie viele Trainingsdatenpunkte der Knoten in jeder Kategorie enthält: Beispielsweise enthält der Knoten rechts unten 0 Exemplare von *Iris setosa*, 1 von *Iris versicolor* und 45 von *Iris virginica*. Schließlich misst das Attribut *gini* die *Gini-Unreinheit* eines Knotens: Ein Knoten gilt als »rein« (gini=0), wenn sämtliche enthaltenen Datenpunkte der gleichen Kategorie angehören. Da beispielsweise der linke Knoten bei depth 1 ausschließlich Instanzen von *Iris setosa* enthält, ist er rein und besitzt eine Gini-Unreinheit von 0. Formel 6-1 stellt dar, wie der Trainingsalgorithmus die Gini-Unreinheit  $G_i$  für den  $i$ . Knoten berechnet. Beispielsweise ist die Gini-Unreinheit für den linken Knoten bei depth 2 gleich  $1 - (0/54)^2 - (49/54)^2 - (5/54)^2 \approx 0,168$ .

Formel 6-1: Gini-Unreinheit

$$G_i = 1 - \sum_{k=1}^n p_{i,k}^2$$

In dieser Gleichung gilt:

- $G_i$  ist die Gini-Unreinheit für den  $i$ . Knoten.
- $p_{i,k}$  ist der Anteil von Instanzen der Kategorie  $k$  an den Datenpunkten im Knoten  $i$ .



Scikit-Learn verwendet den CART-Algorithmus, der ausschließlich *Binärbäume* erzeugt: Innere Knoten haben stets zwei Kinder (d. h., die Fragen lassen sich nur mit Ja oder Nein beantworten). Es existieren aber auch andere Algorithmen wie ID3, die Entscheidungsbäume erzeugen, deren Knoten mehr als zwei Kinder haben können.

Abbildung 6-2 zeigt die Entscheidungsgrenzen dieses Entscheidungsbaums. Die dicke vertikale Linie steht für die Entscheidungsgrenze der Wurzel (depth 0): petal length = 2,45 cm. Da das Gebiet auf der linken Seite rein ist (ausschließlich *Iris setosa*), lässt es sich nicht weiter aufteilen. Das Gebiet auf der rechten Seite ist dagegen unrein, daher teilt es der Knoten bei depth 1 auf der Höhe petal width = 1,75 cm auf (als gestrichelte Linie dargestellt). Da `max_depth` auf 2 gesetzt wurde, endet der Entscheidungsbaum an dieser Stelle. Wenn Sie allerdings `max_depth` auf 3 setzen, würden die zwei Knoten bei depth 2 jeweils eine zusätzliche Entscheidungsgrenze produzieren (hier durch die beiden vertikalen gepunkteten Linien dargestellt).

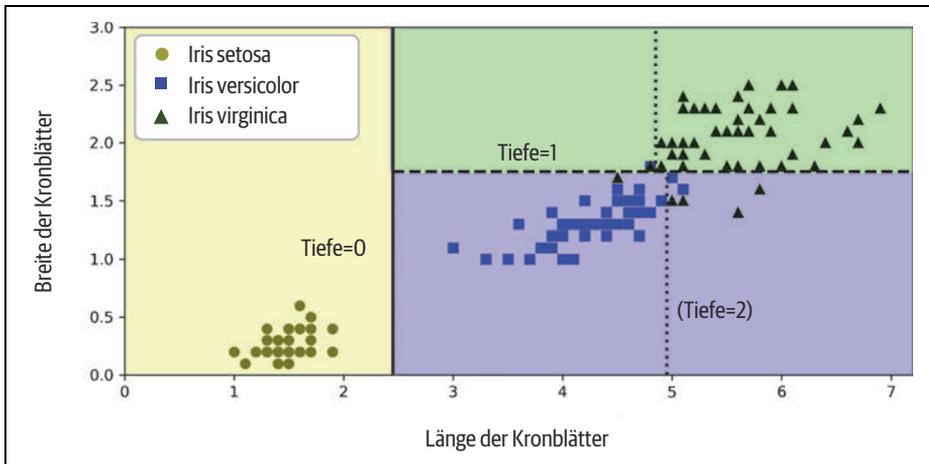


Abbildung 6-2: Entscheidungsgrenzen in einem Entscheidungsbaum



Die Baumstruktur mit all den Informationen, die in Abbildung 6-1 zu sehen sind, steht im Attribut `tree_` des Klassifikators zur Verfügung. Für Details geben Sie `help(tree_clf.tree_)` ein, und im Notebook dieses Kapitels (<https://homl.info/colab3>) finden Sie ein Beispiel.

## Interpretation von Modellen: White Box im Vergleich zu Black Box

Wie Sie sehen, sind Entscheidungsbäume recht einfach nachzuvollziehen und ihre Entscheidungen leicht interpretierbar. Solche Modelle werden oft als *White-Box-Modelle* bezeichnet. Im Gegensatz dazu gehören, wie wir noch sehen werden, Random Forests und neuronale Netze im Allgemeinen zu den *Black-Box-Modellen*. Sie treffen zwar ausgezeichnete Vorhersagen, und Sie können die dazu durchgeführten Berechnungen leicht nachprüfen, dennoch ist es in der Regel schwierig, in wenigen Worten zu erklären, warum die Vorhersagen in einer bestimmten Weise getroffen wurden. Wenn beispielsweise ein neuronales Netzwerk behauptet, dass eine bestimmte Person auf einem Bild zu sehen ist, kann man nur schwer erkennen, worauf sich diese Vorhersage stützt: Hat das Modell die Augen der Person erkannt? Den Mund? Die Nase? Die Schuhe? Oder gar das Sofa, auf dem die Person saß? Umgekehrt geben uns Entscheidungsbäume einfache, klare Regeln zur Klassifikation, die notfalls sogar von Hand durchgeführt werden können (z. B. bei der Klassifikation von Blüten). Das Feld des *interpretierbaren ML* versucht, ML-Systeme zu schaffen, die ihre Entscheidungen für Menschen verständlich erläutern können. Das ist in vielen Bereichen wichtig – zum Beispiel um sicherzustellen, dass das System keine unfairen Entscheidungen trifft.

## Schätzen von Wahrscheinlichkeiten für Kategorien

Ein Entscheidungsbaum kann auch die Wahrscheinlichkeit für die Zugehörigkeit eines Datenpunkts zu einer Kategorie  $k$  abschätzen: Zuerst schreitet das Verfahren den Baum ab, um das Blatt für diesen Datenpunkt zu finden, und gibt dann den Anteil der Trainingsdatenpunkte der Kategorie  $k$  in diesem Knoten zurück. Nehmen wir an, Sie hätten eine Blüte mit 5 cm langen und 1,5 cm breiten Kronblättern entdeckt. Der dazu passende Knoten im Entscheidungsbaum ist der Knoten bei depth 2 auf der linken Seite, daher gibt der Entscheidungsbaum die folgenden Wahrscheinlichkeiten aus: 0% für *Iris setosa* (0/54), 90,7% für *Iris versicolor* (49/54) und 9,3% für *Iris virginica* (5/54). Wenn Sie nach einer Vorhersage der Kategorie fragen, erhalten Sie natürlich *Iris versicolor* (Kategorie 1), da diese die höchste Wahrscheinlichkeit hat. Überprüfen wir dies:

```
>>> tree_clf.predict_proba([[5, 1.5]]).round(3)
array([[0.    , 0.907, 0.093]])
>>> tree_clf.predict([[5, 1.5]])
array([1])
```

Perfekt! Beachten Sie, dass die geschätzten Wahrscheinlichkeiten überall in der unteren rechten Ecke von Abbildung 6-2 die gleichen sind – sogar wenn die Kronblätter beispielsweise 6 cm lang und 1,5 cm breit wären (obwohl es sich dann höchstwahrscheinlich um eine *Iris virginica* handeln würde).

# Der CART-Trainingsalgorithmus

Scikit-Learn verwendet den Algorithmus *Classification and Regression Tree* (CART-Algorithmus), um Entscheidungsbäume zu trainieren (oder »anzubauen«). Er folgt einer sehr einfachen Grundidee: Der Algorithmus teilt die Trainingsdaten zunächst anhand eines Merkmals  $k$  und eines Schwellenwerts  $t_k$  in zwei Untermengen auf (z. B. »petal length  $\leq 2,45$  cm«). Wie werden  $k$  und  $t_k$  ausgewählt? Der Algorithmus sucht nach dem Paar  $(k, t_k)$ , das die reinsten (nach deren Größe gewichteten) Untermengen hervorbringt. Dabei versucht der Algorithmus, die Kostenfunktion in Formel 6-2 zu minimieren.

Formel 6-2: Kostenfunktion des CART-Algorithmus zur Klassifikation

$$J(k, t_k) = \frac{m_{\text{links}}}{m} G_{\text{links}} + \frac{m_{\text{rechts}}}{m} G_{\text{rechts}}$$

wobei  $\begin{cases} G_{\text{links/rechts}} & \text{die Unreinheit der linken/rechten Untermenge misst,} \\ m_{\text{links/rechts}} & \text{die Anzahl Datenpunkte in der linken/rechten Untermenge ist.} \end{cases}$

Sobald der Trainingsdatensatz erfolgreich zweigeteilt wurde, werden die Untermengen nach dem gleichen Verfahren weiter aufgeteilt. Dies setzt sich rekursiv fort, bis die (durch den Hyperparameter `max_depth` angegebene) maximale Tiefe erreicht ist oder keine Aufteilung gefunden werden kann, die die Unreinheit weiter reduziert. Andere Hyperparameter (die wir in Kürze erklären) steuern weitere Abbruchbedingungen: `min_samples_split`, `min_samples_leaf`, `min_weight_fraction_leaf` und `max_leaf_nodes`.



Wie Sie sehen, gehört der CART-Algorithmus zu den Greedy-Algorithmen: Er sucht gierig nach einer optimalen Aufteilung auf der höchsten möglichen Ebene und setzt diesen Vorgang auf jeder Stufe fort. Es wird nicht geprüft, ob eine Aufteilung einige Schritte weiter zur höchsten möglichen Unreinheit führt. Ein Greedy-Algorithmus liefert oft eine recht gute Lösung, diese muss aber nicht die bestmögliche sein.

Leider gehört das Finden des optimalen Baums zu den *NP-vollständigen* Problemen:<sup>1</sup> Es erfordert eine Zeit von  $O(\exp(m))$ , wodurch das Problem selbst für eher kleine Datensätze praktisch unlösbar wird. Daher müssen wir uns mit einer »annehmbaren guten« Lösung zufriedengeben, wenn wir Entscheidungsbäume trainieren.

---

1 P ist die Menge der in polynomieller Zeit lösbaren Probleme. NP ist die Menge der Probleme, deren Lösungen sich in polynomieller Zeit überprüfen lassen. Ein NP-schweres Problem ist ein Problem, das sich in polynomieller Zeit auf ein bekanntes NP-schweres Problem reduzieren lässt. Ein NP-vollständiges Problem ist sowohl NP als auch NP-schwer. Ob  $P = NP$  gilt, ist eine wichtige mathematische Frage. Falls  $P \neq NP$  gilt (was wahrscheinlich erscheint), kann es für kein NP-vollständiges Problem jemals einen polynomiellen Algorithmus geben (außer vielleicht auf einem Quantencomputer).

# Komplexität der Berechnung

Das Treffen von Vorhersagen erfordert das Abschreiten eines Entscheidungsbaums von der Wurzel zu einem Blatt. Entscheidungsbäume sind halbwegs ausbalanciert, daher erfordert das Abschreiten etwa  $O(\log_2(m))$  Knoten, wobei  $\log_2(m)$  der binäre Logarithmus von  $m$  ist, der  $\log_2(m) = \log(m) / \log(2)$  entspricht. Da bei jedem Knoten nur ein Merkmal geprüft werden muss, ist die Komplexität der Vorhersage lediglich  $O(\log_2(m))$ , egal wie viele Merkmale es gibt. Die Vorhersagen sind also auch für große Trainingsdatensätze sehr schnell.

Allerdings vergleicht der Trainingsalgorithmus bei jedem Knoten sämtliche Merkmale (falls `max_features` gesetzt ist, auch weniger) aller Datenpunkte miteinander. Dies führt zu einer Komplexität von  $O(n \times m \log(m))$  beim Trainieren.

## Gini-Unreinheit oder Entropie?

Standardmäßig nutzt der `DecisionTreeClassifier` die Gini-Unreinheit, Sie können aber stattdessen auch die *Entropie* als Maß für die Unreinheit auswählen, indem Sie den Hyperparameter `criterion` auf "entropy" setzen. Der aus der Thermodynamik stammende Begriff Entropie beschreibt Unordnung auf molekularer Ebene: Die Entropie nähert sich null, wenn Moleküle unbeweglich und wohlgeordnet sind. Dieses Konzept hat sich später in andere Fachgebiete ausgebreitet, darunter Shannons *Informationstheorie*, wo sie den durchschnittlichen Informationsgehalt einer Nachricht misst, wie wir in Kapitel 4 gesehen haben. Wenn alle Nachrichten identisch sind, beträgt die Entropie null. Im Machine Learning wird die Entropie oft als Maß für die Unreinheit eingesetzt: Die Entropie einer Menge ist null, wenn sie nur aus Datenpunkten einer Kategorie besteht. Formel 6-3 zeigt die Definition der Entropie des  $i$ . Knotens. Beispielsweise beträgt die Entropie für den linken Knoten bei depth 2 in Abbildung 6-1 genau  $-\frac{49}{54} \log_2\left(\frac{49}{54}\right) - \frac{5}{54} \log_2\left(\frac{5}{54}\right) \approx 0,445$ .

Formel 6-3: Entropie

$$H_i = - \sum_{\substack{k=1 \\ p_{i,k} \neq 0}}^n p_{i,k} \log_2(p_{i,k})$$

Sollten Sie also die Gini-Unreinheit oder die Entropie verwenden? Tatsächlich macht es meist keinen großen Unterschied: Beide ergeben ähnliche Bäume. Die Gini-Unreinheit lässt sich ein wenig schneller berechnen und eignet sich daher als Standardwert. Wenn beide Maße voneinander abweichen, neigt die Gini-Unreinheit dazu, die häufigste Kategorie in einem eigenen Ast des Baums abzusondern, wohingegen die Entropie etwas ausbalanciertere Bäume erzeugt.<sup>2</sup>

---

2 Details finden Sie in einer interessanten Analyse (<https://homl.info/19>) von Sebastian Raschka.

# Hyperparameter zur Regularisierung

Entscheidungsbäume treffen sehr wenige Annahmen über die Trainingsdaten (im Gegensatz zu beispielsweise linearen Modellen, die annehmen, dass sich die Daten linear verhalten). Sich selbst überlassen, passt sich die Struktur des Baums sehr genau an die Trainingsdaten an und führt höchstwahrscheinlich zu Overfitting. Solch ein Modell wird auch als *parameterfreies Modell* bezeichnet, nicht weil es keine Parameter gäbe (meist gibt es viele), sondern weil die Anzahl der Parameter nicht vor dem Trainieren festgelegt wird. Daher ist es dem Modell überlassen, sich eng an die Daten anzupassen. Im Gegensatz dazu besitzt ein *parametrisches Modell* wie etwa ein lineares Modell eine im Voraus festgelegte Anzahl Parameter. Es verfügt daher über eine begrenzte Anzahl Freiheitsgrade, wodurch es weniger zu Overfitting neigt (aber dafür ein höheres Risiko für Underfitting besteht).

Um ein Overfitting der Trainingsdaten zu vermeiden, müssen Sie die Freiheitsgrade eines Entscheidungsbaums beim Trainieren einschränken. Wie Sie inzwischen wissen, nennt man dies Regularisierung. Die Hyperparameter zur Regularisierung hängen vom verwendeten Algorithmus ab. Im Allgemeinen können Sie aber mindestens die maximale Tiefe des Entscheidungsbaums begrenzen. In Scikit-Learn lässt sich dies über den Hyperparameter `max_depth` erreichen. Die Voreinstellung ist `None`, was für eine unbegrenzte Tiefe steht. Ein Reduzieren von `max_depth` reguliert das Modell und verringert damit das Risiko einer Überanpassung.

Die Klasse `DecisionTreeClassifier` bietet einige weitere Parameter, die die Form des Entscheidungsbaums in ähnlicher Weise einschränken:

- `max_features`  
Maximale Anzahl der Merkmale, die beim Aufteilen eines Knoten berücksichtigt werden.
- `max_leaf_nodes`  
Maximale Anzahl Blätter.
- `min_samples_split`  
Minimale Anzahl an Datenpunkten, die ein Knoten aufweisen muss, damit er aufgeteilt werden kann.
- `min_samples_leaf`  
Minimale Anzahl an Datenpunkten, die ein Blatt haben muss.
- `min_weight_fraction_leaf`  
Wie `min_samples_leaf`, aber als Anteil der gesamten gewichteten Datenpunkte.

Ein Erhöhen der `min_*`-Hyperparameter und ein Senken der `max_*`-Hyperparameter reguliert das Modell.



Andere Algorithmen trainieren Entscheidungsbäume zunächst ohne Einschränkungen, entfernen aber anschließend überflüssige Knoten (*Pruning*). Ein Knoten wird als überflüssig angesehen, wenn seine Kinder ausschließlich Blätter sind und der durch ihn erbrachte Zugewinn an Reinheit nicht *statistisch signifikant* ist. Standardisierte statistische Tests wie der  $\chi^2$ -Test schätzen die Wahrscheinlichkeit ab, dass eine Verbesserung rein zufällig erfolgt (dies bezeichnet man als die *Nullhypothese*). Wenn diese Wahrscheinlichkeit, genannt *p-Wert*, höher als ein eingestellter Schwellenwert ist (typischerweise 5%, durch einen Hyperparameter steuerbar), dann wird der Knoten als überflüssig erachtet, und seine Kinder werden gelöscht. Dieses Pruning wird fortgesetzt, bis alle überflüssigen Knoten entfernt worden sind.

Testen wir die Regularisierung auf dem *moons*-Datensatz, der in Kapitel 5 vorgestellt wurde. Wir werden einen Entscheidungsbaum ohne Regularisierung trainieren und einen anderen mit `min_samples_leaf=5`. Dies ist der Code, Abbildung 6-3 zeigt die Entscheidungsgrenzen jedes Baums:

```
from sklearn.datasets import make_moons

X_moons, y_moons = make_moons(n_samples=150, noise=0.2, random_state=42)

tree_clf1 = DecisionTreeClassifier(random_state=42)
tree_clf2 = DecisionTreeClassifier(min_samples_leaf=5, random_state=42)
tree_clf1.fit(X_moons, y_moons)
tree_clf2.fit(X_moons, y_moons)
```

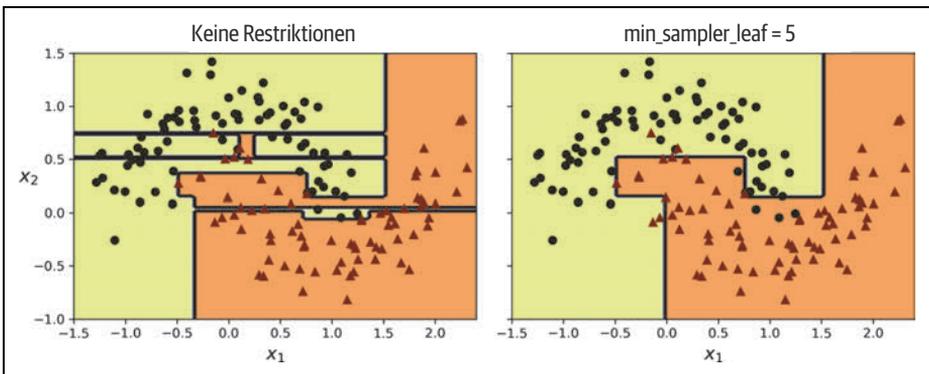


Abbildung 6-3: Entscheidungsgrenzen eines nicht regularisierten (links) und eines regularisierten Baums (rechts)

Das unregularisierte Modell links ist ganz klar overfittet, während das regularisierte Modell rechts wahrscheinlich besser verallgemeinern wird. Wir können das überprüfen, indem wir beide Bäume mit einem Testdatensatz evaluieren, der mit einem anderen Zufalls-Seed erzeugt wurde:

```
>>> X_moons_test, y_moons_test = make_moons(n_samples=1000, noise=0.2,
...                                          random_state=43)
```

```

...
>>> tree_clf1.score(X_moons_test, y_moons_test)
0.898
>>> tree_clf2.score(X_moons_test, y_moons_test)
0.92

```

Tatsächlich hat der zweite Baum eine bessere Genauigkeit für den Testdatensatz.

## Regression

Entscheidungsbäume können auch Regressionsaufgaben bewältigen. Erstellen wir mit der Klasse `DecisionTreeRegressor` aus Scikit-Learn einen Regressionsbaum und trainieren wir diesen auf einem verrauschten quadratischen Datensatz mit `max_depth=2`:

```

import numpy as np
from sklearn.tree import DecisionTreeRegressor

np.random.seed(42)
X_quad = np.random.rand(200, 1) - 0.5 # ein einzelnes, zufälliges Eingabemerkmakl
y_quad = X_quad ** 2 + 0.025 * np.random.randn(200, 1)

tree_reg = DecisionTreeRegressor(max_depth=2, random_state=42)
tree_reg.fit(X_quad, y_quad)

```

Der dabei erhaltene Baum ist in Abbildung 6-4 dargestellt.

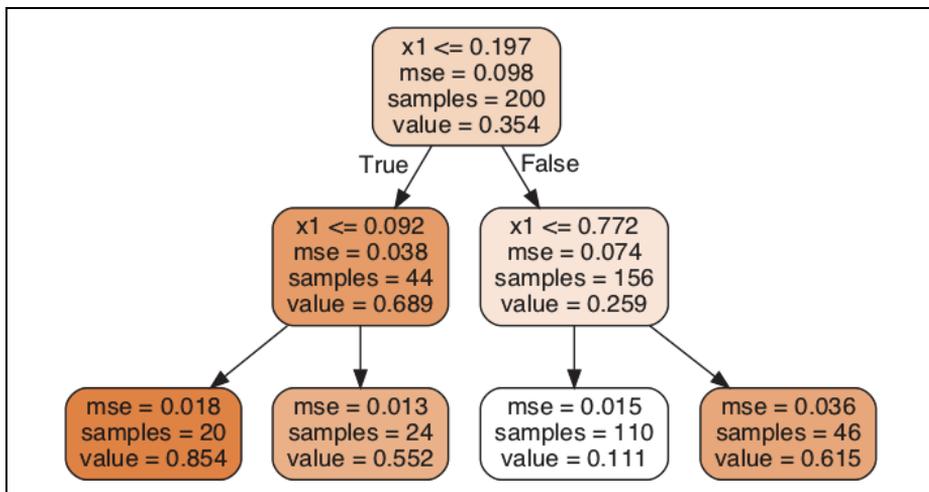


Abbildung 6-4: Ein Entscheidungsbaum zur Regression

Dieser Baum sieht dem zuvor zur Klassifikation erstellten Baum sehr ähnlich. Der Hauptunterschied ist, dass er anstelle einer Kategorie in jedem Knoten einen Wert vorhersagt. Wenn Sie beispielsweise eine Vorhersage für einen neuen Datenpunkt bei  $x_1 = 0,2$  treffen möchten, fragt die Wurzel, ob  $x_1 \leq 0,197$ . Da das nicht der Fall

ist, wechselt der Algorithmus zum rechten Kindknoten, der fragt, ob  $x_1 \leq 0,772$ . Das stimmt, daher geht der Algorithmus zum linken Kindknoten über. Das ist ein Blatt, und es sagt  $\text{value}=0.111$  vorher. Diese Vorhersage ist der durchschnittliche Zielwert der 110 Trainingsdatenpunkte in diesem Blatt. Sie führt zu einem mittleren quadratischen Fehler (MSE) von 0,015 in den 110 Datenpunkten.

Die Vorhersagen dieses Modells sind auf der linken Seite von Abbildung 6-5 dargestellt. Wenn Sie  $\text{max\_depth}=3$  setzen, erhalten Sie die auf der rechten Seite dargestellten Vorhersagen. Beachten Sie, dass der vorhergesagte Wert in jedem Abschnitt dem durchschnittlichen Zielwert der Datenpunkte in diesem Abschnitt entspricht. Der Algorithmus teilt jeden Abschnitt so auf, dass möglichst viele Trainingsdatenpunkte so nah wie möglich am vorhergesagten Wert liegen.

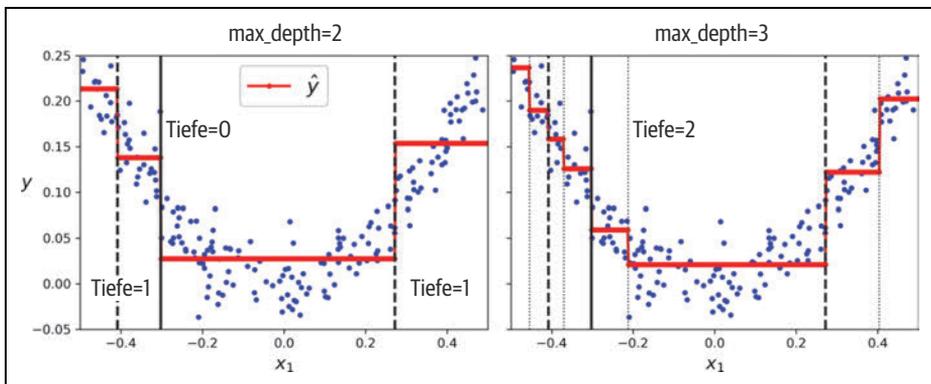


Abbildung 6-5: Vorhersagen zweier Entscheidungsbäume zur Regression

Der CART-Algorithmus funktioniert wie weiter oben beschrieben. Der einzige Unterschied besteht darin, dass der Trainingsdatensatz so aufgeteilt wird, dass der MSE statt der Unreinheit minimiert wird. Formel 6-4 zeigt die vom Algorithmus minimierte Kostenfunktion.

Formel 6-4: CART-Kostenfunktion für die Regression

$$J(k, t_k) = \frac{m_{\text{links}}}{m} \text{MSE}_{\text{links}} + \frac{m_{\text{rechts}}}{m} \text{MSE}_{\text{rechts}} \quad \text{mit} \quad \begin{cases} \text{MSE}_{\text{Knoten}} = \frac{\sum_{i \in \text{Knoten}} (\hat{y}_{\text{Knoten}} - y^{(i)})^2}{m_{\text{Knoten}}} \\ \hat{y}_{\text{Knoten}} = \frac{\sum_{i \in \text{Knoten}} y^{(i)}}{m_{\text{Knoten}}} \end{cases}$$

Wie bei Klassifikationsaufgaben sind auch Entscheidungsbäume für Regressionsaufgaben anfällig für Overfitting. Ohne jegliche Regularisierung (z. B. mit den voreingestellten Hyperparametern) erhalten Sie die Vorhersagen auf der linken Seite von Abbildung 6-6. Dies ist offensichtlich ein schwerwiegendes Overfitting der Trainingsdaten. Durch Setzen von  $\text{min\_samples\_leaf}=10$  erhalten Sie ein weitaus sinnvollerer Modell, das Sie auf der rechten Seite von Abbildung 6-6 sehen.

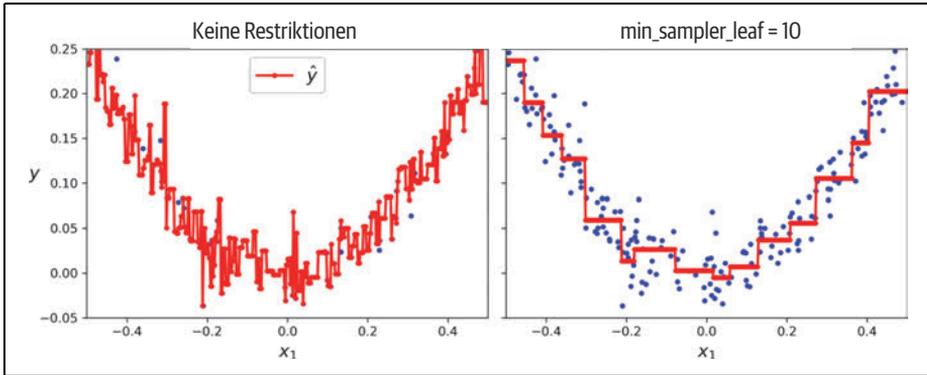


Abbildung 6-6: Vorhersagen eines unregulierten Regressionsbaums (links) und eines regulierten Baums (rechts)

## Empfindlichkeit für die Ausrichtung der Achsen

Sie sind hoffentlich mittlerweile davon überzeugt, dass sehr vieles für Entscheidungsbäume spricht: Sie sind recht einfach zu verstehen und zu interpretieren, leicht anwendbar und mächtig. Sie haben allerdings auch einige Schwachstellen: Erstens haben Sie womöglich bereits festgestellt, dass Entscheidungsbäume orthogonale Entscheidungsgrenzen lieben (alle Unterteilungen stehen im rechten Winkel zu einer Achse). Dadurch reagieren sie empfindlich auf die Orientierung der Daten. Als Beispiel zeigt Abbildung 6-7 einen einfachen linear separierbaren Datensatz: Auf der linken Seite kann der Entscheidungsbaum diesen einfach unterteilen. Auf der rechten Seite wirkt die Entscheidungsgrenze nach einer Drehung des Datensatzes um  $45^\circ$  unnötig verschachtelt. Obwohl beide Entscheidungsbäume die Trainingsdaten perfekt abbilden, verallgemeinert das Modell auf der rechten Seite vermutlich nicht besonders gut.

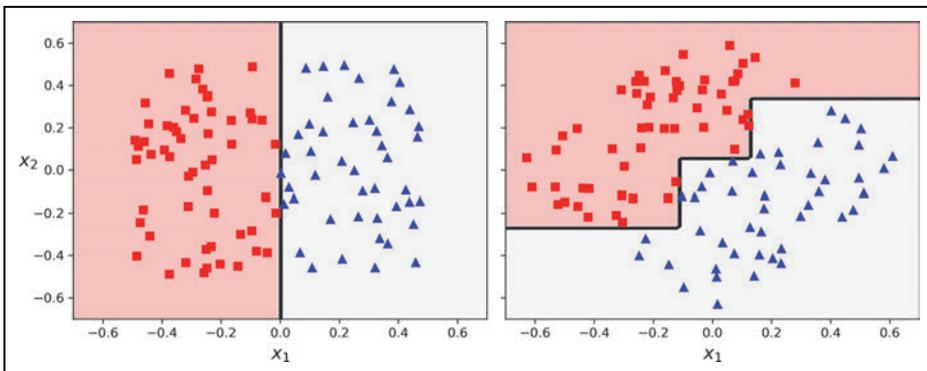


Abbildung 6-7: Empfindlichkeit für eine Rotation der Trainingsdaten

Eine Möglichkeit, dieses Problem anzugehen, ist das Skalieren der Daten und das anschließende Anwenden einer Hauptkomponentenzerlegung. Wir werden auf diese genauer in Kapitel 8 eingehen, aber hier müssen Sie nur wissen, dass sie die Daten so rotiert, dass sich Bäume danach oft (aber nicht immer) leichter erstellen lassen.

Bauen wir eine kleine Pipeline, die die Daten skaliert und mittels PCA rotiert, um danach einen `DecisionTreeClassifier` damit zu trainieren. In Abbildung 6-8 sehen Sie die Entscheidungsgrenzen dieses Baums: Die Rotation ermöglicht es, den Datensatz ziemlich gut mit nur einem Merkmal  $z_1$  abzubilden, bei dem es sich um eine lineare Funktion der ursprünglichen Länge und Breite des Kronblatts handelt. Dies ist der Code:

```
from sklearn.decomposition import PCA
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import StandardScaler

pca_pipeline = make_pipeline(StandardScaler(), PCA())
X_iris_rotated = pca_pipeline.fit_transform(X_iris)
tree_clf_pca = DecisionTreeClassifier(max_depth=2, random_state=42)
tree_clf_pca.fit(X_iris_rotated, y_iris)
```

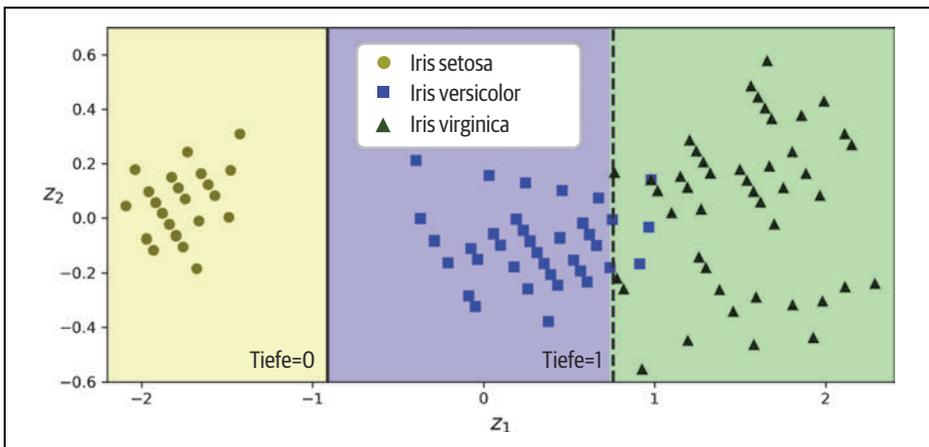


Abbildung 6-8: Die Entscheidungsgrenzen eines Baums mit dem skalierten und per PCA rotierten Iris-Datensatz

## Entscheidungsbäume haben eine größere Varianz

Allgemeiner gesprochen, ist das Hauptproblem bei Entscheidungsbäumen, dass sie eine ziemlich große Varianz besitzen: Kleine Änderungen an den Hyperparametern oder den Daten können zu sehr unterschiedlichen Modellen führen. Genauer gesagt: Da der von Scikit-Learn verwendete Trainingsalgorithmus stochastisch arbeitet – die bei jedem Knoten zu evaluierenden Merkmale werden zufällig ausgewählt –,

können Sie sogar mit den gleichen Trainingsdaten sehr unterschiedliche Modelle erhalten, wie es in Abbildung 6-9 zu sehen ist (es sei denn, Sie setzen den Hyperparameter `random_state`). Dieser Entscheidungsbaum unterscheidet sich stark vom vorherigen (Abbildung 6-2).

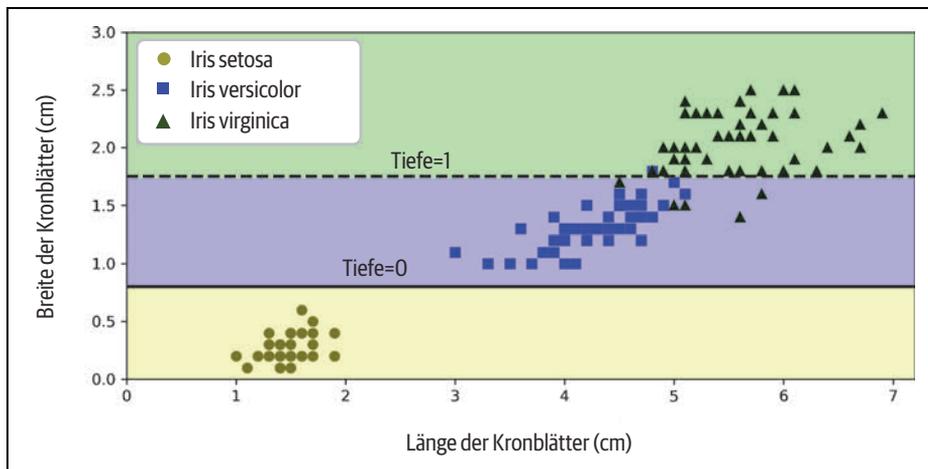


Abbildung 6-9: Ein erneutes Trainieren des gleichen Modells mit den gleichen Daten kann ein sehr anderes Modell erzeugen.

Glücklicherweise ist es durch das Mitteln der Vorhersagen über viele Bäume möglich, die Varianz deutlich zu verringern. Solch ein *Ensemble* aus Bäumen, das *Random Forest* genannt wird, ist eine der leistungsfähigsten Arten von heutzutage verfügbaren Modellen, wie Sie im nächsten Kapitel sehen werden.

## Übungen

1. Was ist die ungefähre Tiefe eines mit 1 Million Datenpunkten (ohne Restriktionen) trainierten Entscheidungsbaums?
2. Ist die Gini-Unreinheit eines Knotens im Allgemeinen geringer oder größer als die seines Elternteils? Ist sie *im Allgemeinen* kleiner/größer oder *immer* kleiner/größer?
3. Sollte man versuchen, `max_depth` zu senken, wenn ein Entscheidungsbaum einen Trainingsdatensatz overfittet?
4. Sollte man versuchen, die Eingabemerkmale zu skalieren, wenn ein Entscheidungsbaum die Trainingsdaten underfittet?
5. Wenn es eine Stunde dauert, einen Entscheidungsbaum mit 1 Million Datenpunkten zu trainieren, wie lange etwa wird das Trainieren eines weiteren Baums mit 10 Millionen Datenpunkten dauern? Hinweis: Denken Sie an die Berechnungskomplexität des CART-Algorithmus.

6. Wenn es eine Stunde dauert, einen Entscheidungsbaum mit einem gegebenen Trainingsdatensatz zu trainieren, wie lange dauert es dann ungefähr, wenn sich die Anzahl der Merkmale verdoppelt?
7. Trainieren und optimieren Sie einen Entscheidungsbaum für den Datensatz *moons* anhand dieser Schritte.
  - a. Erzeugen Sie einen *moons*-Datensatz mit `make_moons(n_samples=10000, noise=0.4)`.
  - b. Teilen Sie ihn mit `train_test_split()` in einen Trainings- und einen Testdatensatz auf.
  - c. Verwenden Sie die Gittersuche mit Kreuzvalidierung (mithilfe der Klasse `GridSearchCV`), um gute Einstellungen für die Hyperparameter eines `DecisionTreeClassifier` zu finden. Hinweis: Probieren Sie unterschiedliche Werte für `max_leaf_nodes`.
  - d. Trainieren Sie den Baum mit den vollständigen Trainingsdaten und bestimmen Sie die Qualität Ihres Modells auf den Testdaten. Sie sollten eine Genauigkeit zwischen 85% und 87% erhalten.
8. Züchten Sie einen Wald mit diesen Schritten.
  - a. Erzeugen Sie im Anschluss an die vorherige Aufgabe 1.000 Untermengen Ihres Trainingsdatensatzes mit jeweils 100 zufällig ausgewählten Datenpunkten. Hinweis: Sie können dazu die Klasse `ShuffleSplit` aus `Scikit-Learn` verwenden.
  - b. Trainieren Sie auf jedem der Teildatensätze einen Entscheidungsbaum mit den besten oben gefundenen Hyperparametern. Werten Sie diese 1.000 Entscheidungsbäume auf den Testdaten aus. Da sie mit kleineren Datensätzen trainiert wurden, schneiden diese Bäume mit nur ca. 80% voraussichtlich schlechter ab als der erste Entscheidungsbaum.
  - c. Nun kommt der Zaubertrick. Generieren Sie für jeden Testdatenpunkt die Vorhersagen der 1.000 Entscheidungsbäume und heben Sie ausschließlich die häufigste Vorhersage auf (Sie können dazu die Funktion `mode()` aus `SciPy` verwenden). Damit erhalten Sie *Mehrheitsvorhersagen* für Ihren Testdatensatz.
  - d. Werten Sie diese Vorhersagen mit dem Testdatensatz aus: Sie sollten eine etwas höhere Genauigkeit als beim ersten Modell erhalten (etwa 0,5 bis 1,5% höher). Herzlichen Glückwunsch, Sie haben soeben einen `Random-Forest-Klassifikator` trainiert!

Die deutschsprachigen Lösungen zu diesen Übungen finden Sie unter <https://dpunkt.de/produkt/praxiseinstieg-machine-learning-mit-scikit-learn-keras-und-tensorflow-2/> auf der Registerkarte *Zusatzmaterial*, die englischsprachigen Lösungen am Ende des Notebooks zu diesem Kapitel unter <https://homl.info/colab3>.