



Prof. Kristian Kersting  
Steven Lang  
Felix Friedrich

Sommersemester 2022  
Übungsblatt 5

## Benötigte Dateien

Alle benötigten Datensätze und Skriptvorlagen finden Sie in unserem Moodle-Kurs:

<https://moodle.informatik.tu-darmstadt.de/course/view.php?id=1058>

## 5.1 Bagging - Zufallswälder

Ein Zufallswald (engl. Random Forrest) ist ein Meta-Schätzer, der eine Reihe von Entscheidungsbaum-Klassifikatoren auf verschiedene Unterstichproben des Datensatzes anpasst und die Mittelwertbildung verwendet, um die Vorhersagegenauigkeit zu verbessern und Überanpassung (engl. Overfitting) zu kontrollieren. Dies wird als Bagging-Methode bezeichnet. Bei Ensemble-Algorithmen bilden Bagging-Methoden eine Klasse von Algorithmen, die mehrere Instanzen eines Black-Box-Schätzers auf zufälligen Untermengen des ursprünglichen Trainingssatzes aufbauen und dann ihre individuellen Vorhersagen zu einer endgültigen Vorhersage aggregieren. Diese Methoden werden als eine Möglichkeit verwendet, die Varianz eines Basisschätzers (z.B. eines Entscheidungsbaums) zu reduzieren, indem man die Randomisierung in sein Konstruktionsverfahren einführt und dann ein Ensemble daraus macht. In vielen Fällen stellen Bagging-Methoden eine sehr einfache Möglichkeit zur Verbesserung gegenüber einem einzelnen Modell dar, ohne dass der zugrundeliegende Basisalgorithmus angepasst werden muss. Da sie eine Möglichkeit zur Verringerung der Überanpassung bieten, funktionieren Bagging-Methoden am besten mit starken und komplexen Modellen (z.B. tiefen Entscheidungsbäumen), im Gegensatz zu Boosting-Methoden, die normalerweise am besten mit schwachen Modellen (z.B. flachen Entscheidungsbäumen) funktionieren.

### a) Entscheidungsbaum vs. Zufallswald (diskret)

Wir möchten in einer einfachen Demonstration die Performanz von Entscheidungsbäumen mit Zufallswäldern vergleichen. Wir werden beide Modelle anhand eines Satzes von Klassifizierungs- und Regressionsdatensätzen kreuzvalidieren.

Implementieren Sie die Funktion `cross_val_dt_rt`, welche eine 10-fache Kreuzvalidierung für einen `DecisionTreeClassifier` mit Standardparametern und `RandomForestClassifier` mit 20 Schätzern durchführt. Sie können dabei die Funktion `sklearn.model_selection.cross_validate` verwenden. Abschließend sollen die Genauigkeitsscores zurückgegeben werden.

```
1 def cross_val_dt_rt(load_fun) -> (dict, dict):
2     """Conducts a 10-fold cross validation for a decision tree and random forest and
3     returns their accuracy scores given an sklearn discrete dataset loader function."""
4     # Load data
5     X, y = load_fun(return_X_y=True)
6     # Cross-validate Decision Tree and Random Forest
7     scores_dt = cross_validate(
8         estimator=DecisionTreeClassifier(),
9         X=X,
10        y=y,
11        scoring="accuracy",
12        cv=10,
13        return_train_score=True,
14        n_jobs=-1,
15    )
16    scores_rf = cross_validate(
17        estimator=RandomForestClassifier(n_estimators=20),
18        X=X,
19        y=y,
20        scoring="accuracy",
```

```

21         cv=10,
22         return_train_score=True,
23         n_jobs=-1,
24     )
25     return scores_dt, scores_rf

```

### b) Entscheidungsbaum vs. Zufallswald (kontinuierlich)

Es wird deutlich, dass ein Random Forest einen Entscheidungsbaum in Testvorhersagegenauigkeit übertrifft. Wir können dies für Datensätze mit kontinuierlichen Zielvariablen wiederholen.

Führen Sie eine 10-fache Kreuzvalidierung in der Funktion `cross_val_dt_rf_continuous` für kontinuierliche Zielvariablen durch. Es sollen hierbei erneut 20 Schätzer verwendet werden und die mittlere quadratische Abweichung als Metrikscore zurückgegeben werden.

```

1 def cross_val_dt_rf_continuous(load_fun) -> (dict, dict):
2     """Conducts a 10-fold cross validation for a decision tree and random forest and
3     returns their mse scores given an sklearn discrete dataset loader function."""
4     # Load data
5     X, y = load_fun(return_X_y=True)
6     # Cross-validate Decision Tree and Random Forest
7     scores_dt = cross_validate(
8         estimator=DecisionTreeRegressor(),
9         X=X,
10        y=y,
11        scoring=make_scorer(mean_squared_error),
12        cv=10,
13        return_train_score=True,
14        n_jobs=-1,
15    )
16    scores_rf = cross_validate(
17        estimator=RandomForestRegressor(n_estimators=20),
18        X=X,
19        y=y,
20        scoring=make_scorer(mean_squared_error),
21        cv=10,
22        return_train_score=True,
23        n_jobs=-1,
24    )
25    return scores_dt, scores_rf

```

## 5.2 Zufallswälder - Hyperparameter

Was wir beobachten, ist, dass der Random Forest nicht nur den Entscheidungsbaum übertrifft, sondern auch die Trainingsdaten nicht so sehr überanpasst, wie der Entscheidungsbaum. Für jeden Datensatz passt sich der einzelne Baum vollständig an die Trainingsdaten an und erreicht einen mittleren quadratischen Wurzelfehler von 0.0, während der Zufallswald nicht vollständig mit den Trainingsdaten übereinstimmt. Auf der anderen Seite erzielt der Zufallswald in jedem Fall ein besseres Testergebnis. Dies deutet darauf hin, dass Zufallswälder bessere Verallgemeinerungsfähigkeiten haben als Entscheidungsbäume.

### a) Anzahl der Bäume

Der wichtigste Parameter eines Zufallswaldes ist seine Ensemblegröße, d.h. wie viele Entscheidungsbäume verwendet werden, um eine Mehrheitsabstimmung für die Entscheidung durchzuführen. Beginnen wir mit dem Zifferndatensatz und zeigen wir den Einfluss der Ensemblegröße auf die Klassifikationsgenauigkeit.

Um den Einfluss der Anzahl der Bäume zu sehen, werden wir jeden Wert im Intervall  $n \in [1, \dots, 40]$  kreuzvalidieren und ein `RandomForestClassifier`-Modell mit  $n$ -Entscheidungsbäumen erstellen. Evaluieren Sie die Anzahl der Schätzer in Funktion `evaluate_n_estimators` mittels einer 10-fachen Kreuzvalidierung und geben sie die mittlere Trainings- und Testgenauigkeit zurück.

```

1 def evaluate_n_estimators(X: np.ndarray, y: np.ndarray, n: int) -> (float, float):
2     """Run 10 fold cross-validation of the model for a given number of trees and returns the
3     mean train and test score."""
4     model = RandomForestClassifier(n_estimators=n)

```

```

5     scores = cross_validate(
6         estimator=model, X=X, y=y, scoring="accuracy", cv=10, return_train_score=True, n_jobs=-1,
7     )
8     return np.mean(scores["train_score"]), np.mean(scores["test_score"])

```

Was können Sie bei zunehmender Anzahl von Bäumen im Ensemble feststellen?

*Wir können beobachten, dass mit zunehmender Anzahl von Bäumen im Ensemble auch die Trainings- und Testgenauigkeit zunimmt. Da die Erhöhung der Anzahl der Bäume zu einer Verringerung der Verzerrung und Varianz führt, stellen wir auch fest, dass das Modell mit zunehmender Komplexität nicht an die Trainingsdaten überangepasst wird.*

#### b) Maximal Baumtiefe

Ein weiterer interessanter Aspekt, den es zu betrachten gilt, ist die Komplexität der Entscheidungsbäume. Wir können den Einfluss der Komplexität der Lerner auf die Leistung des Ensembles untersuchen, indem wir das obige Experiment wiederholen, aber statt mit einer unterschiedlichen Anzahl von Bäumen zu evaluieren, werden wir diesmal die maximale Baumtiefe für jeden Baum mit einer festen Anzahl von Bäumen ändern.

Führen Sie eine 10-fache Kreuzvalidierung für gegebene Baumtiefe mit 20 Schätzern in der Funktion `evaluate_depth` durch. Was fällt ihnen bei größer werdender Baumtiefe auf?

```

1 def evaluate_depth(X: np.ndarray, y: np.ndarray, depth: int) -> (float, float):
2     """Run 10 fold cross-validation of the model for a given tree depth and returns the
3     mean train and test score."""
4     model = RandomForestClassifier(max_depth=depth, n_estimators=20)
5     scores = cross_validate(
6         estimator=model, X=X, y=y, scoring="accuracy", cv=10, return_train_score=True, n_jobs=-1,
7     )
8     return np.mean(scores["train_score"]), np.mean(scores["test_score"])

```

*Es wird deutlich, dass bei Zufallswäldern das Ensemble umso stärker wird, je stärker die Entscheidungsbäume sind, ohne dass die Verallgemeinerung verloren geht. Dies ist eine wichtige Beobachtung, da bei anderen Algorithmen, die auf einem Satz von Basislernern basieren, schwache Lerner bevorzugt werden können, oder leistungsschwache Lerner ebenso wie starke Lerner, wie z.B. im Fall von Boosting.*

**c) Anzahl der Merkmale**

Um die Korrelation zu reduzieren und Zufälligkeit in das Ensemble zu induzieren, wählt der Zufallswald-Algorithmus eine zufällige Teilmenge von  $k$ -Merkmalen aus allen verfügbaren Eingabemerkmalen für jeden internen Entscheidungsbaum aus. Die Verwendung einer geringeren Anzahl von Eingabemerkmalen verringert die Ähnlichkeit zwischen den einzelnen Bäumen, führt aber auch zu weniger komplexen und daher schwächeren Bäumen. Andererseits wird durch die Erhöhung der Anzahl der Merkmale jeder Baum leistungsstärker, aber auch die Korrelation zwischen den Bäumen erhöht (die bei Verwendung *aller* Merkmale maximiert wird).

Evaluieren Sie den Einfluss dieses Parameters, indem Sie eine 10-fache Kreuzvalidierung mit 20 Schätzern für alle mögliche Anzahl von Merkmalen im Zifferndatensatz in der Funktion `evaluate_features` durchführen.

```
1 def evaluate_features(X: np.ndarray, y: np.ndarray, n_features: int) -> (float, float):
2     """Run 10 fold cross-validation of the model for a given number of features per tree and returns
3     the
4     mean train and test score."""
5     model = RandomForestClassifier(max_features=n_features, n_estimators=20)
6     scores = cross_validate(
7         estimator=model, X=X, y=y, scoring="accuracy", cv=10, return_train_score=True, n_jobs=-1,
8     )
9     return np.mean(scores["train_score"]), np.mean(scores["test_score"])
```

Welche Faustregel wird in der Literatur oft verwendet, um die Anzahl der Merkmale festzulegen? Ist diese Regel auch in diesem Fall sinnvoll?

Die Ergebnisse bestätigen, dass ein Wert um den gebildeten Schätzwert von  $k = \sqrt{n_{\text{features}}} = \sqrt{64} = 8$ , den man in der wissenschaftlichen Literatur bezüglich eines guten Wertes für die Anzahl der Merkmale finden kann, zu einer guten Wahl von  $k$  führt.