



Diese Übung wird am 10.06.2021 um 13:30 Uhr besprochen und nicht bewertet.

Benötigte Dateien

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

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

6.1 Adaboost - Entscheidungsgrenzen

Das Kernprinzip von AdaBoost besteht darin, eine Sequenz schwacher Lerner (d. h. Modelle, die nur geringfügig besser sind als zufällige Vermutungen, wie z. B. flache Entscheidungsbäume) auf wiederholt modifizierte Versionen der Daten anzupassen. Die von allen schwachen Lernenden erhaltenen Vorhersagen werden dann durch eine gewichtete Mehrheitsentscheidung (oder Summe) kombiniert, um die endgültige Vorhersage zu erstellen. Die Datenmodifikationen bei jeder sogenannten Boosting-Iteration bestehen aus der Anwendung von Gewichten w_1, w_2, \dots, w_N zu jeder der N Trainingsproben. Zu Beginn werden diese Gewichte alle auf $w_i = 1/N$ gesetzt, so dass im ersten Schritt lediglich ein schwacher Lerner auf den Originaldaten trainiert wird. Für jede aufeinanderfolgende Iteration werden die Stichprobengewichte individuell modifiziert, und der Lernalgorithmus wird erneut auf die neu gewichteten Daten angewendet. In einem Schritt werden die Gewichte der Trainingsbeispiele, die durch das im vorherigen Schritt induzierte verstärkte Modell falsch vorhergesagt wurden, erhöht, während die Gewichte für diejenigen, die korrekt vorhergesagt wurden, verringert werden. Im Laufe der Iterationen erhalten schwer vorhersagbare Beispiele einen immer stärkeren Einfluss. Jeder nachfolgende schwache Lerner ist dadurch gezwungen, sich auf die Beispiele zu konzentrieren, die von den vorhergehenden in der Sequenz übersehen werden [1].

a) Datensatzvisualisierung

Beginnen wir mit einem künstlichen Datensatz, der aus zwei konzentrischen Kreisen besteht, die zwei Klassen repräsentieren:

Visualisieren Sie den Datensatz in der Methode `plot_dataset`. Die Punkt der ersten und zweiten Klasse sollen dabei farblich unterschieden werden und Trainings- und Testpunkte verschieden markiert werden.

b) Basislerner

Als Basislerner für das Boosting-Modell werden wir einen einfachen Entscheidungsbaum der Tiefe 1 (wegen seiner visuellen Flachheit auch Entscheidungsstumpf genannt) wählen. Um ein Gefühl für dessen Einfachheit zu bekommen, können wir den Stumpf auf die Daten selbst anpassen und die Entscheidungsgrenze visualisieren.

Visualisieren Sie die Entscheidungsgrenze unseres Basislernalgorithmus in der Methode `plot_decision_boundary_stump`.

c) Algorithmus Beschreibung

Es wird deutlich, dass ein Entscheidungsstumpf allein nicht in der Lage ist, die beiden Klassen voneinander zu trennen. Da der Stumpf die Tiefe 1 hat, kann er die Daten nur auf der Basis eines **einzigen** Merkmals trennen, das in diesem Fall etwa $X_1 \approx 0.5$ beträgt. Jeder Datenpunkt mit $X_1 \leq 0.5$ wird als blaue Klasse vorhergesagt, während alles mit $X_2 > 0.55$ als rote Klasse vorhergesagt wird.

Daher werden wir nun den AdaBoost-Algorithmus verwenden, um diesen einfachen Basislerner wie oben beschrieben zu verbessern.

Beschreiben Sie in eigenen Worten das Verfahren des Adaboost-Algorithmus anhand Abbildung 1.

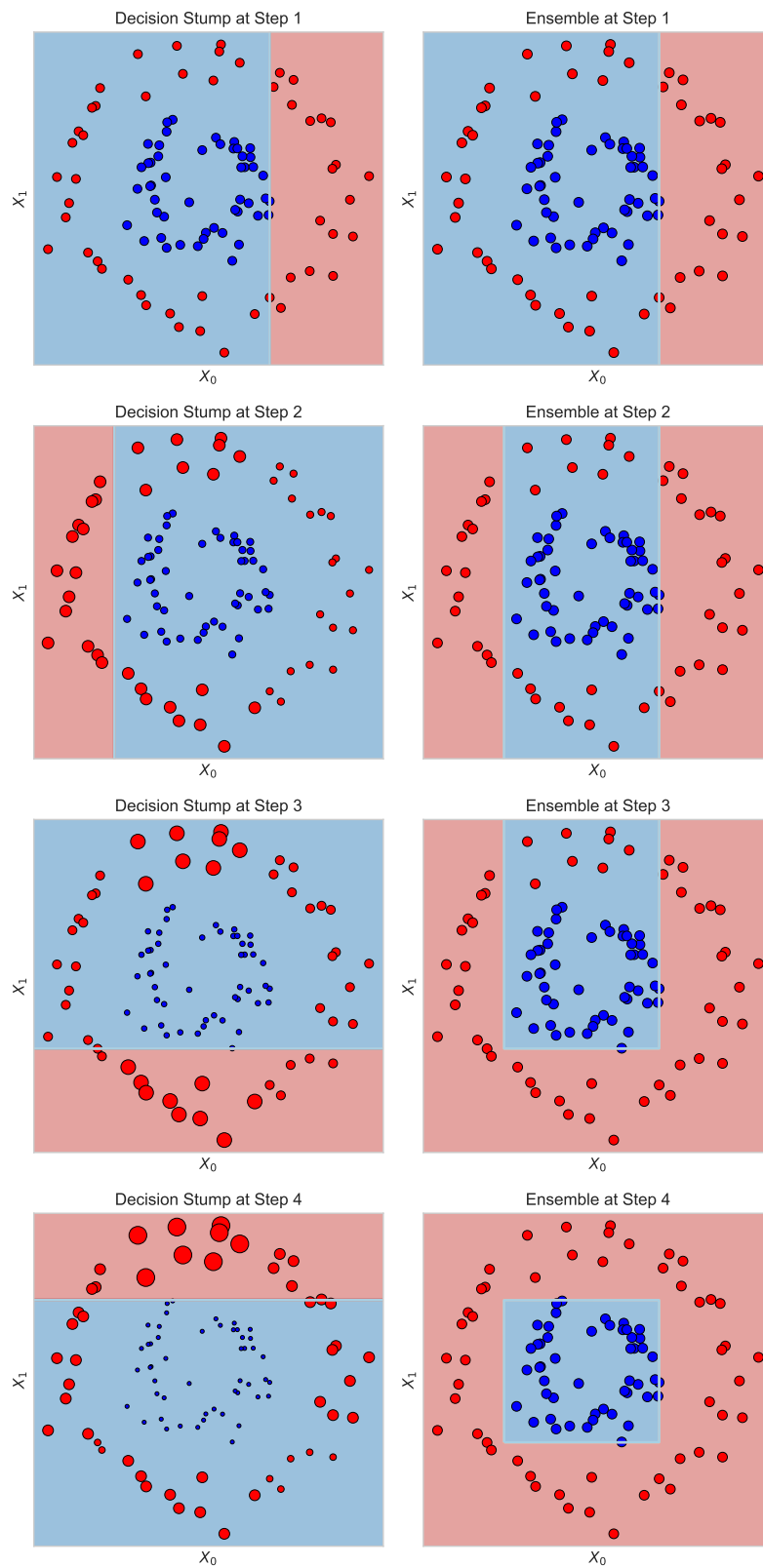


Abbildung 1: Iteration von Adaboost mit Entscheidungsstümpfen.

6.2 Adaboost - Klassifikation

Schauen wir uns ein etwas komplexeres Beispiel an: Das Klassifizieren von Ziffern. Der Datensatz `sklearn.datasets.load_digits` enthält 1797 verschiedene 8×8 Pixel-Bilder von Ziffern.

a) **Datensatzvisualisierung**

Visualisieren Sie die ersten fünf Trainingsbeispiele mit ihrem zugehörigen Label in der Methode `plot_digit_dataset`.

b) **Iterationen / Anzahl Klassifier**

Wir werden das AdaBoost-Modell erneut mit einem Entscheidungstumpf als Basis-Lerner für 300 Iterationen lernen. Für jede Iteration wird ein zusätzlicher Entscheidungstumpf eingesetzt.

Sklearn ermöglicht es uns über wiederholte Anwendung der Methode `AdaBoostClassifier.staged_predict` stufenweise Vorhersagen zu treffen. Dies erzeugt Vorhersagen nach jeder Boosting-Iteration und erlaubt es uns, die Genauigkeitskurve über die Trainingsiterationen zu visualisieren.

Fitten Sie ein `AdaBoostClassifier` über die Methode `get_acc_for_estimators` und geben Sie das Modell und die zugehörigen Trainings- und Testgenauigkeiten für jede Iteration zurück.

Wie verhält sich die Performanz mit ansteigender Anzahl von Iterationen?

c) **Lernrate**

Es ist wichtig zu beachten, dass die Steilheit und die tatsächliche Verbesserung auch von der Lernrate abhängt, d.h. der Rate, mit der jeder zusätzliche Basislernende zur Vorhersage beiträgt. Um zu sehen, wie viel Einfluss die Lernrate hat, können wir das Diagramm aus b) für verschiedene Lernraten visualisieren: Fitten Sie `AdaBoostClassifier` über Ihre zuvor implementierte Funktion `get_acc_for_estimators` für die Lernraten 2^i wobei $i \in \{-3, -2, -1, 0, 1\}$ und geben Sie die Trainingsgenauigkeiten für jede Iteration zurück. Welche Schlussfolgerung können Sie aus dem Plot zum Vergleich der Lernraten ziehen?

6.3 Gradient Boosting

Gradient Tree Boosting oder gradientenverstärkte Regressions Bäume (GBRT) ist eine Verallgemeinerung des Boosting zu willkürlich differenzierbaren Verlust-Funktionen. GBRT ist ein präzises und effektives Verfahren, das sowohl für Regressions- als auch für Klassifikationsprobleme verwendet werden kann. GBRT Modelle werden in einer Vielzahl von Bereichen eingesetzt, darunter Web-Suchrangfolge und Ökologie.

Die Vorteile des GBRT sind:

- Natürliche Handhabung von Daten gemischten Typs (= heterogene Merkmale)
- Gute Vorhersagekraft
- Robustheit gegenüber Ausreißern im Zielvariablenraum (durch robuste Verlustfunktionen)

Die Nachteile des GBRT sind:

- Skalierbarkeit, aufgrund der sequentiellen Natur des Boostings kann es kaum parallelisiert werden.

Algorithmus

1. Trainieren Sie einen Entscheidungstumpf (Entscheidungsbaum mit einer Tiefe von 1) T_1 auf den Daten. Dies führt zu unserem ersten Modell $M_1 = T_1$. Die Schwachstellen sind die Residuen $r = y_{true} - y_{pred}$.
2. Trainieren Sie einen weiteren Entscheidungstumpf T_2 auf den Residuen. Dies führt zu unserem zweiten Modell $M_2 = M_1 + \gamma_2 T_2$, wobei γ die Schrittweite ist. Als Regularisierung können wir eine Lernrate $\nu \in (0, 1)$ (manchmal als Schrumpfung bezeichnet) einführen und erhalten $M_2 = M_1 + \nu \gamma_2 T_2$.
3. Wiederholen Sie 2. und trainieren Sie weitere Entscheidungstümpfe auf den Residuen von M_2 , dann M_3 , bis eine gute Anpassung an die Zielvariable erreicht ist

a) Regression

Beginnen wir mit dem Lernen einer einfachen Sinuskurve.

Mit dem Modul `GradientBoostingRegressor` von `sklearn` können wir ein GradientBoosting-Regressionsmodell an die erzeugte Sinuskurve anpassen. Wir werden die Residuen für insgesamt 100 Schritte anpassen und die Funktion der kleinsten Quadrate als Fehlerfunktion verwenden.

Wenden Sie einen Gradient Boosting Regressor in der Funktion `fit_gradient_boosting_regressor` an und geben Sie das gelernte Model und die Trainingsfehler für jede Iteration zurück.

b) Konzept des Algorithmus

Um den iterativen Fortschritt des Modells auf den Trainingsdaten zu visualisieren, ist es möglich, den Fortschritt nach jeder Iteration zu erhalten (s. Abb. 2).

Beschreiben Sie das Konzept des Algorithmus von Gradient Boosting.

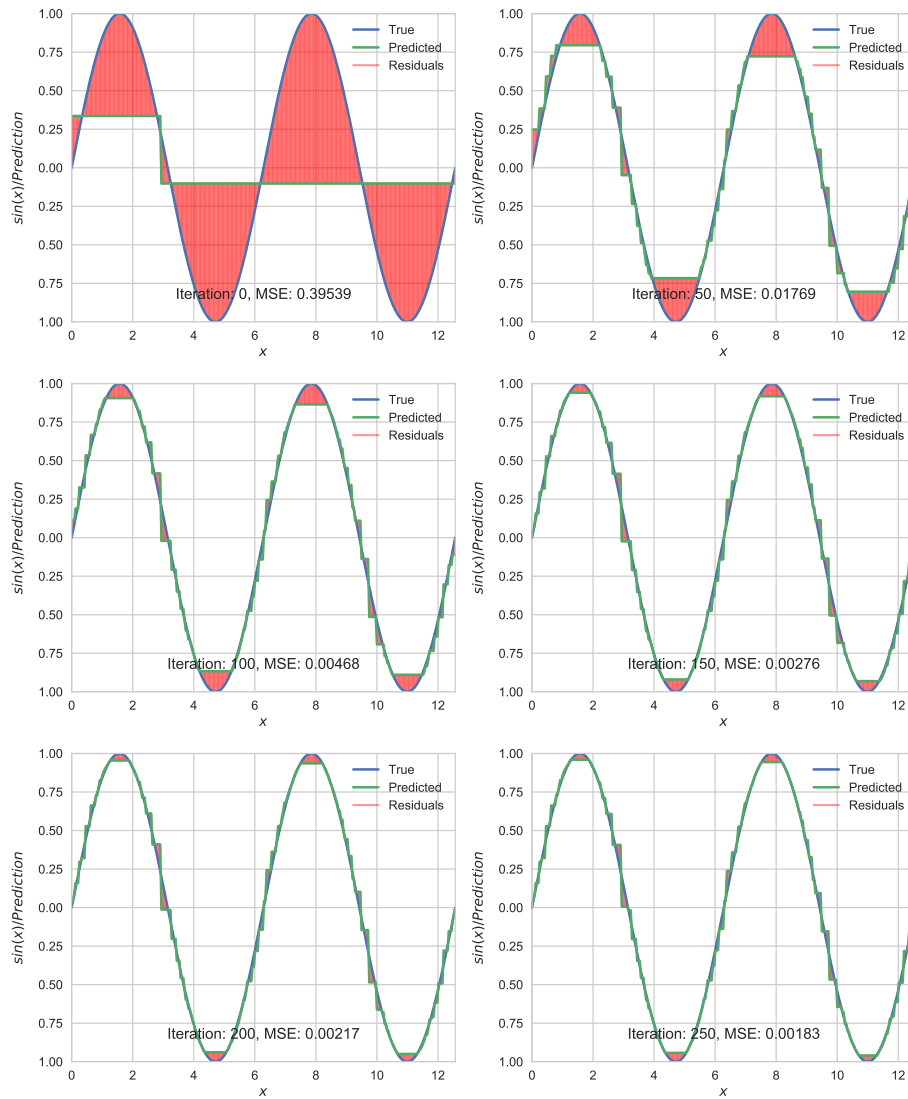


Abbildung 2: Iterationen von Gradient Boost mit Entscheidungsstümpfen.

6.4 Gradient Boosting - Fehlerfunktionen

Die folgenden Fehlerfunktionen werden für Gradient Boosting unterstützt [1] und können in `sklearn` mit dem Parameter 'loss' spezifiziert werden:

Regression

- Kleinste Quadrate (engl. Least Squares) ('ls'): Die natürliche Wahl für die Regression aufgrund ihrer überlegenen rechnerischen Eigenschaften. Das Ausgangsmodell ist durch den Mittelwert der Zielwerte gegeben.
- Geringste absolute Abweichung ('lad'): Eine robuste Verlustfunktion für die Regression. Das Anfangsmodell ist durch den Median der Zielwerte gegeben.
- Huber ('huber'): Eine robuste Verlustfunktion für die Regression: Eine weitere robuste Verlustfunktion, die kleinste Quadrate und geringste absolute Abweichung kombiniert; verwenden Sie Alpha, um die Sensitivität in Bezug auf Ausreißer zu kontrollieren (siehe [F2001] für weitere Details).
- Quantil ('quantil'): Eine Verlustfunktion für die Quantil-Regression. Verwenden Sie $0 < \alpha < 1$, um das Quantil zu spezifizieren. Diese Verlustfunktion kann verwendet werden, um Prädiktionsintervalle zu erstellen.

Klassifikation

- Binomische Abweichung ('deviance'): Die negative Binomial-Log-Wahrscheinlichkeitsverlust-Funktion für binäre Klassifikation (liefert Wahrscheinlichkeitsschätzungen). Das anfängliche Modell ist durch das Log-Wahrscheinlichkeits-Verhältnis gegeben.
- Multinomiale Abweichung ('deviance'): Die negative binomiale logarithmische Wahrscheinlichkeitsverlust-Funktion für binäre Klassifikation (liefert Wahrscheinlichkeitsschätzungen): Die negative multinomiale logarithmische Wahrscheinlichkeitsverlustfunktion für die Mehrklassenklassifikation mit $n_classes$ sich gegenseitig ausschließender Klassen. Sie liefert Wahrscheinlichkeitsschätzungen. Das anfängliche Modell ist durch die vorherige Wahrscheinlichkeit jeder Klasse gegeben. Bei jeder Iteration müssen $n_Klassen$ -Regressionsbäume konstruiert werden, was die GBRT für Datensätze mit einer großen Anzahl von Klassen ziemlich ineffizient macht.
- Exponentieller Verlust ('exponential'): Dieselbe Verlustfunktion wie AdaBoostClassifier. Weniger robust gegenüber falsch vorhergesagten Datenpunkte als 'deviance'; kann nur für binäre Klassifikation verwendet werden.

Regularisierung

[2] schlug eine einfache Regularisierungsstrategie vor, die den Beitrag jedes schwachen Lernalgorithmus um einen Faktor ν skaliert:

$$F_m(x) = F_{m-1}(x) + \nu \gamma_m h_m(x)$$

a) Fehlerfunktionen

Evaluieren Sie die Performanz von Gradient Boosting in der Funktion `evaluate_loss_fn` unter Nutzung der verschiedenen, oben genannten Fehlerfunktionen und geben Sie jeweils die mittleren Fehlerquadrate (MSE) für den Trainings- und Testdatensatz zurück.

b) Lernrate

In folgender Abbildung 3 wird eine Auswertung des Bostoner Datensatzes für Train- und Test-MSE unter Verwendung von Lernraten $lr \in \{2^{-5}, \dots, 2^0\}$ gezeigt. Beurteilen Sie wie sich unterschiedliche Lernraten auswirken.

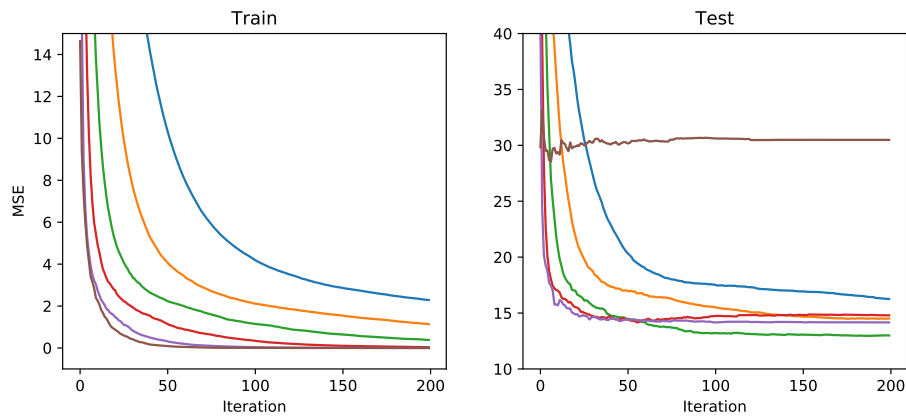


Abbildung 3: Performanz von Gradient Boosting mittels MSE unter Anwendung unterschiedlicher Lernraten auf dem Boston Datensatz.

Literatur

- [1] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.
- [2] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.