

# Data Mining und Maschinelles Lernen

Prof. Kristian Kersting  
Steven Lang  
Felix Friedrich



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

---

Sommersemester 2022  
Übungsblatt 3

---

## Benötigte Dateien

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

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

---

## 3.1 Bewertungsmetriken

---

Die Wahl der richtigen Bewertungsmetrik ist ein wichtiger Schritt beim Einsatz von maschinellen Lernsystemen. Dies wird die Art und Weise prägen, wie Ergebnisse interpretiert und Modelle ausgewählt und verfeinert werden. Jede Metrik hat ihre eigenen Vor- und Nachteile, deren man sich bewusst sein muss.

### a) Klassifikation

In einer binären Klassifikation mit den Klassen  $\{pos, neg\}$  kann ein Modell vier Arten von Vorhersagen machen: **Wahr Positiv (TP)**, **Wahr Negativ (TN)**, **Falsch Positiv (FP)**, **Falsch Negativ (FN)**. Erklären Sie kurz diese Begriffe und bewerten Sie, ob sie zur Bewertung gleichbedeutend sind.

### b) Entscheidungsschwelle

Die meisten Klassifikatoren geben nicht einfach den Klassenwert selbst zurück, sondern eher eine Wahrscheinlichkeit oder Pseudo-Wahrscheinlichkeit, ob die Stichprobe  $x_i$  zur  $pos$ -Klasse gehört. Diese Wahrscheinlichkeit ist ein Wert zwischen 0.0 und 1.0. Es ist daher notwendig, einen Schwellenwert zu wählen, bei dem wir entscheiden, wann eine Probe zur Klasse  $pos$  gehört. Wenn dieser Schwellenwert z.B. bei 0.5 liegt und der Klassifizierer die Wahrscheinlichkeit  $P(pos|x_i) = 0.3$  zuweist, würden wir die  $neg$ -Klasse vorhersagen. Erläutern Sie wie Sie eine solche Entscheidungsschwelle beispielsweise in der Krebsdiagnose einsetzen können.

---

### 3.2 Precision, Recall, Genauigkeit

---

Weiterhin gibt es die Metriken **Precision**, **Recallwert**, **F1-Wertung**, **Genauigkeit** (engl. Precision, Recall, Accuracy, F1 Score). Erläutern Sie deren Bedeutung und wie diese definiert sind:

a) **Precision**

b) **Recall**

c) **F1-Wert**

d) **Genauigkeit**

e) **Beispiel**

Sie arbeiten an einem Modell, das in Bildern Unkraut und Nutzpflanzen voneinander unterscheiden soll. In der Evaluation schnitt das Modell wie folgt ab: Von 135 getesteten Bildern wurden 45 korrekt als Unkraut, 50 korrekt als Nutzpflanze, 25 inkorrekt als Unkraut, und 15 inkorrekt als Nutzpflanze klassifiziert. Berechnen Sie die Precision, Recall, und F1-Werte der Klassifikation. Hinweis: Sie können die Nutzpflanzen als positive Klasse ansehen.

f) **PR-Kurve**

Die Optimierung auf Precision führt in der Regel zu einer Verringerung des Recalls und umgekehrt. Dies zeigt sich in einem PR-Kurvendiagramm wie im folgenden Beispiel, in dem wir einen RandomForestClassifier auf dem Brustkrebs-Datensatz ausführen. Die Datenpunkte für die PR-Kurve erhält man durch Berechnung der Precision und des Recall-Wertes für einen Satz von Schwellenwerten zwischen 0.0 und 1.0.

Implementieren Sie die Funktion `plot_pr_curve`. Sie können dabei die Funktion `sklearn.metrics.precision_recall_curve` verwenden.

g) **ROC-Kurve**

Die ROC-Kurve (engl. Receiver Operating Characteristic) besteht aus der Falsch-Positiv Rate, die gegen die Wahr-Positiv Rate aufgetragen wird:

- **Wahr-Positiv Rate (TPR)** (äquivalent zu Recall):

$$TPR = \frac{TP}{TP + FN}$$

- **Falsch-Positive Rate (FPR)**:

$$FPR = \frac{FP}{FP + TN}$$

Ähnlich wie bei der PR-Kurve werden die TPR- und FPR-Werte für einen Bereich für Schwellenwerte zwischen 0.0 und 1.0 berechnet.

Unter Verwendung der ROC-Kurve ist es möglich, den sogenannten **AUC-Score** (engl. Area under Curve) als alternative Metrik zur Genauigkeit zu berechnen. Der AUC-Score misst die Fläche unter der ROC-Kurve (höher ist besser). Im Gegensatz zum Genauigkeitsscore bietet er eine Messung über alle möglichen Schwellenwerte hinweg.

Implementieren Sie die Funktion `plot_roc_curve` um die ROC-Kurve zu visualisieren. Geben Sie im Plot-Titel den zugehörigen AUC-Wert an. Sie können dabei die Funktionen `sklearn.metrics.roc_curve` und `sklearn.metrics.roc_auc_score` verwenden.

h) **Konfusionsmatrix**

Die Konfusionsmatrix (engl. confusion matrix) zeigt, wie oft eine wahre Klasse mit einer anderen Klasse verwechselt wurde. Die y-Achse stellt die wahren Beschriftungen dar, während die x-Achse die vorhergesagten Beschriftungen zeigt.

Zeigen Sie die Konfusionsmatrix über die Funktion `plot_confusion_matrix` an. Sie können dabei die Funktionen `sklearn.metrics.confusion_matrix` und `seaborn.heatmap` verwenden.

### 3.3 Fehlermetriken der Regression

Für Regressionsaufgaben gibt es mehrere Metriken, wie die Differenz zwischen dem wahren Zielwert  $y_i$  einer Stichprobe  $x_i$  und dem vorhergesagten Zielwert  $\hat{y}_i$  über alle Datenpunkte aggregiert werden kann.

- **Mittlerer absoluter Fehler (engl. mean absolute error (MAE))**

- Behandelt alle Residuen gleich

$$\text{MAE}(y, \hat{y}) = \frac{1}{N} \sum_{i=0}^{N-1} |y_i - \hat{y}_i|$$

- **Mittlerer quadratischer Fehler (MSE) (engl. mean squared error (MSE))**

- Bestraft große Residuen stärker durch quadratischen Effekt

$$\text{MSE}(y, \hat{y}) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2$$

- **Wurzel des mittleren quadratischen Fehlers (engl. root mean squared error (RMSE))**

- Gleich wie MSE, aber in der gleichen Einheit wie die Zielvariable  $y$

$$\text{RMSE}(y, \hat{y}) = \sqrt{\frac{1}{N} \sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2}$$

- **Mittlerer quadratischer logarithmischer Fehler (MSLE)**

- Wird am besten verwendet, wenn das Ziel ein exponentielles Wachstum aufweist
- Bestraft eine zu niedrig vorhergesagte Schätzung höher als eine zu hoch vorhergesagte Schätzung

$$\text{MSLE}(y, \hat{y}) = \frac{1}{N} \sum_{i=0}^{N-1} (\log_e(1 + y_i) - \log_e(1 + \hat{y}_i))^2$$

- **Medianer Absoluter Fehler (engl. median absolute error (MAE))**

- Robust gegenüber Ausreißern

$$\text{MedAE}(y, \hat{y}) = \text{median}(|y_0 - \hat{y}_0|, \dots, |y_{N-1} - \hat{y}_{N-1}|)$$

- **$R^2$  Punktzahl (engl.  $R^2$  Score)**

- Misst die Anpassungsgüte eines Modells
- Optimale Anpassung bei einem Wert von 1
- Schlechter als Mittelwert-Vorhersage, wenn der Wert unter 0 liegt
- Datensatzabhängig aufgrund von Varianz-Term

$$\begin{aligned} R^2(y, \hat{y}) &= 1 - \frac{\sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2}{\sum_{i=0}^{N-1} (y_i - \bar{y})^2} \\ &= 1 - \frac{\text{MSE}(y, \hat{y})}{\text{Var}[y]} \end{aligned}$$

#### a) Metrikevaluierung

Evaluieren Sie die oben beschriebenen Metriken in der Funktion `evaluate_metric`. Wenden Sie dabei zunächst eine 5-fache Kreuzvalidierung mittels `sklearn.model_selection.cross_val_predict` an und werten Sie anschließend die Metrik über die Parameterübergabe von `y_true` und `y_pred` aus.