

Wörterliste ML Typische Probleme: Bias / Variance und Class Imbalance

Underfitting / grosser Bias	Modelle mit niedriger Komplexität sind robuster, aber manchmal zu einfach.
Underfitting – was hilft?	<ul style="list-style-type: none"> - Mehr Trainingsdaten - Komplexeres Modell wählen bspw. Gradient Boosting
Overfitting / grosse Varianz	Modelle mit hoher Komplexität sind genauer, generalisieren aber manchmal nicht/schlecht.
Overfitting – was hilft?	<ul style="list-style-type: none"> - Komplexität reduzieren.. <ul style="list-style-type: none"> o Einfachere Modellart bspw. Logistische Regression o Gleiches Modell aber vereinfachen (regularisieren, prunen) o Verwendung von Validierungsmenge, um die richtigen Werte der entsprechenden Parameter zu finden
Class Imbalance	<ul style="list-style-type: none"> - Problem: «Nadeln im Heuhaufen finden», interessante Klasse hat nur wenig Beispiele (Bswp. Marketing, oft viel mehr nein werte als ja) <p>Was nun?: Statt nur „Accuracy“ (Wenn du sehr viele Beispiele von der Mehrheitsklasse hast, und nur ganz wenige von der Minderheitsklasse, dann kann dein Modell faul sein und trotzdem eine hohe Accuracy haben – obwohl es null nützlich ist.) solltest du:</p> <ol style="list-style-type: none"> 1. Kostenbasierte Bewertung machen, wenn du weißt: Was kostet ein falsch-positives vs. falsch-negatives Ergebnis? 2. Oder nutze Metriken wie: <ul style="list-style-type: none"> o Precision = Wie viele der als positiv erkannten sind wirklich positiv? o Recall = Wie viele der echten positiven wurden gefunden? o F1-Score = Kombo aus beiden o AUC (Area under Curve) = Gesamtgüte der Klassifikation <p>→ Diese Metriken schauen gezielt auf die Minderheitsklasse („die Nadeln“), nicht auf das Ganze.</p>
Oversampling	<p>Oversampling: Die kleine Klasse wird vergrössert</p> <ul style="list-style-type: none"> • Wie?: Vorhandene Beispiele werden kopiert oder neue künstliche Daten erstellen • Gut wenn?: wenig Daten vorhanden sind, keine wichtigen Daten verloren werden sollen

Undersampling	<p>Undersampling: Die grosse Klasse wird verkleinert</p> <ul style="list-style-type: none"> • Wie?: Die Mehrheit der Daten wird weggeworfen, damit beide Klassen gleich gross sind • Gut wenn?: Sehr viele Daten vorhanden sind, Verlust von Informationen nicht schlimm ist
Over-/Undersampling – Konsequenzen	<ul style="list-style-type: none"> - Das Modell hat mehr «Mut» die Minderheitsklasse vorherzusagen. Effekte für die Metriken: <ul style="list-style-type: none"> ○ Accuracy = sinkt ○ AUC = steigt ○ F1 = steigt ○ Kosten = sinken ○ Confusion Matrix = mehr TP/FP
Trainingsmenge	<ul style="list-style-type: none"> - Wird verwendet, um das Modell zu trainieren - Das Modell „lernt“ hier aus den Daten - Alles ist erlaubt, was die Modellleistung verbessert: <ul style="list-style-type: none"> ○ Re-Sampling (z. B. Daten duplizieren) ○ Hinzufügen ähnlicher Daten aus externen Quellen - Beispiel: Seltene Klassen wie „Premium-Kunden“ künstlich häufiger machen.
Validierungsmenge	<ul style="list-style-type: none"> - Wird genutzt, um das Modell während des Trainings zu testen - Dient zur Auswahl und Abstimmung von Parametern - Muss realitätsnah und unverändert bleiben - Beispiel: Verteilung von Alter, Geschlecht oder Kaufverhalten muss wie im echten Anwendungsfall sein
Testmenge	<ul style="list-style-type: none"> - Wird nach dem Training verwendet, zu endgültigen Bewertung des Modells - Zeigt, wie gut das Modell auf völlig neue, echte Daten reagiert - Muss ebenfalls unverändert und repräsentativ sein - Beispiel: Daten echter Kunden, die das Modell noch nie gesehen hat
Re-Sampling	<ul style="list-style-type: none"> - Technik zur Änderung der Datenverteilung im Training - Ziel: Klassen gleichmäßig oder fair vertreten
Verteilung von Attributen	<ul style="list-style-type: none"> - Bedeutet: Wir oft bestimmte Werte in den Daten vorkommen - Muss in Validierungs- und Testdaten bleiben