

Interpretierbarkeit

Target Leakage	<p>Das Modell schummelt, weil es beim Training Infos bekommt, die es beim echten Einsatz nicht zur Verfügung hat. Das Modell wirkt im Training super, aber versagt dann im echten Einsatz.</p>
Argumente für Interpretierbarkeit	<p>Warum ist es wichtig, dass Modelle verständlich und nachvollziehbar sind.</p> <ol style="list-style-type: none"> 1. Debugging <ol style="list-style-type: none"> a) Data Leakage erkennen <ul style="list-style-type: none"> - das bedeutet, dass das Modell Informationen verwendet, die es beim echten Einsatz gar nicht haben darf - Interpretierbarkeit hilft, so etwas zu erkennen b) Umgang mit Ausnahmefällen <ul style="list-style-type: none"> - Manche Situationen im echten Leben sind ungewöhnlich oder selten - Bspw. Ein autonomes Auto erkennt Fahrradfahrer. Das Modell achtet nur auf zwei sichtbare Räder. Problem: Fahrräder mit Gepäcktaschen, die die Räder verdecken, werden nicht erkannt - Verständliche Modelle helfen, solche Schwächen zu entdecken c) Unerwünschten Bias erkennen <ul style="list-style-type: none"> - Bspw. Ein Modell bevorzugt männliche Bewerber, weil es aus verzerrten Daten gelernt hat - Nur wenn wir das Modell verstehen, können wir Diskriminierung verhindern 2. Vertrauen <ul style="list-style-type: none"> - Menschen müssen die Modellentscheidungen nachvollziehen können, bevor sie auf ihnen basierend entscheiden, da sie verantwortlich sind 3. Wissenschaftliche Erkenntnis <ul style="list-style-type: none"> - Durch verständliche Modelle können Forscher und Experten neue Muster oder Zusammenhänge entdecken, die vorher nicht bekannt waren
Arten von Interpretierbarkeit	<p>Intrinsisch(modell kann man anschauen) <-> post-hoc(modell kann man nicht anschauen)</p> <p>Modellspezifisch(funktioniert nur für ein Modell) <-> modell-agnostisch(funktioniert für jedes Modell)</p> <p>Lokal(erklärt eine Instanz) <-> global(erklärt das ganze Modell)</p>
Intrinsisch interpretierbare Modelle (Modelle von Natur aus leicht verständlich sind)	<p>Lineare Modelle (Lineare Regression, Naive Bayes etc.):</p> <p>Vorteile:</p> <ul style="list-style-type: none"> - gut erforscht - gut für numerische Merkmale geeignet - die Gewichte zeigen direkt, wie stark jedes Merkmal die Vorhersage beeinflusst

	<p>Nachteile:</p> <ul style="list-style-type: none"> - Wechselwirkungen zwischen Merkmalen werden nicht erfasst - Kann keine komplexen Zusammenhänge erkennen <p>Logikbasierte Modelle (Entscheidungsbäume):</p> <p>Vorteile:</p> <ul style="list-style-type: none"> - Sehr einfach zu verstehen - Erkennt Interaktionen zwischen Merkmalen <p>Nachteile:</p> <ul style="list-style-type: none"> - Numerische Werte werden immer diskretisiert - Bäume können instabil sein und teilweise schwer zu interpretieren
<p>Permutation Feature Importance</p>	<p>Methode, um zu bestimmten, wie wichtig ein einzelnes Merkmal für ein Modell ist.</p> <p>Idee: Wie stark steigt der Fehler des Modells, wenn man das Attribut «kaputtmacht», also die Werte zufällig mischt.</p> <p>→ wenn ein Feature wichtig ist, dann wird das Modell viel schlechter</p> <p>Vorteile:</p> <ul style="list-style-type: none"> - Wenn ein Feature z. B. nur zusammen mit einem anderen wichtig ist, wird das bei der Permutation trotzdem erkannt, weil die Beziehung dabei auch zerstört wird <p>Nachteile:</p> <ul style="list-style-type: none"> - Benötigt Zugang zu den Labels der Daten - Wenn man zwei stark zusammenhängende Features hat, dann scheint keines von beiden wichtig, weil das Modell die Info doppelt bekommt, das macht die Interpretation nicht immer logisch nachvollziehbar
<p>Partial Dependence Plots</p>	<p>Verstehen, wie ein einzelnes Merkmal die Modellvorhersage beeinflusst, im Durchschnitt über alle Datenpunkte</p> <p>Beispiel:</p> <p>Wenn du wissen willst, wie sich „Alter = 50“ auf die Vorhersage auswirkt, dann:</p> <ul style="list-style-type: none"> - nimm alle Datenpunkte, - setze deren Alter auf 50, - rechne für jeden die Vorhersage neu aus, - und bilde den Durchschnitt dieser Vorhersagen → das ist der PDP-Wert bei „50“. <p>Vorteile:</p> <ul style="list-style-type: none"> - intuitiv verständlich <p>Nachteile:</p> <ul style="list-style-type: none"> - funktioniert nur gut für 1 oder 2 Features gleichzeitig - ignoriert Korrelationen

Shapely-Werte (lokal)	<p>Man stellt sich vor, dass die Merkmale einer Instanz nacheinander in zufälliger Reihenfolge ins Modell «hineingehen»</p> <ul style="list-style-type: none"> - Jedes Feature trägt etwas zur Gesamtvorhersage bei - Der Shapley-Wert eines Merkmals ist: → Wie viel ändert sich die Vorhersage im Schnitt, wenn dieses Feature dem Team beitritt? <p>Beispiel:</p> <ul style="list-style-type: none"> - Wenn Alter „als letztes reinkommt“ – wie stark ändert sich dann die Vorhersage? - Man wiederholt das mit allen möglichen Reihenfolgen und mittelt <p>Vorteile:</p> <ul style="list-style-type: none"> - Mathematisch fundiert - Kann vollständige Erklärungen liefern <p>Nachteile:</p> <ul style="list-style-type: none"> - Rechenintensiv - Nicht sofort intuitiv
SHAP Summary Plot(global)	<ul style="list-style-type: none"> - Jeder Punkt ist ein Shapley- Wert für eine Instanz und ein Attribut. - Attribute sind nach Wichtigkeit absteigend sortiert, Attributwerte farblich kodiert (rot = grosse Werte). - Werte auf der x-Achse entsprechen der Vorhersage(wahrscheinlichkeit)
Individual Conditional Expecatiton	<ul style="list-style-type: none"> - Zeigt eine Linie pro Instanz, aus der man jeweils sieht, wie sich die Vorhersage für diese Instanz ändert, wenn der Attributwert (x- Achse) sich ändert. <p>Vorteile:</p> <ul style="list-style-type: none"> - Zeigt mehr Details (Verteilung) als PDPs - Ziemlich intuitiv <p>Nachteile:</p> <ul style="list-style-type: none"> - Nur 1 Feature darstellbar - Bildet keine Interaktionen/Korrelationen ab