



Meilenstein 3

Verlustprävention an Selbstbedienungskassen im Einzelhandel

*Durchgeführt durch die
Retail Data Mining GmbH*



HEUTE IN BERUFE DER ZUKUNFT:
DER KI-FLÜSTERER

Themen für heute:

1. Vorbemerkung
2. Statische Regeln
3. Klassifikationsmodell
4. Regressionsmodell
5. Gesamtmodell
6. Bewertung und Empfehlung
7. Ausblick

1. Vorbemerkung

Ziele des Meilensteins (1)

- **Ausgangsbasis:** bereinigte & aggregierte Transaktionsdaten
- **Ziel:** Entwicklung **praxistauglicher Modelle** zur Betrugserkennung an SBK

- Berücksichtigung betriebswirtschaftlicher **Bewertungsfunktion**

	Tatsächlich FRAUD	Tatsächlich NORMAL
Vorhersage FRAUD	+5 (TP)	-10 (FP)
Vorhersage NORMAL	-Schaden (FN)	0 (TN)

- Werteverlust reduzieren <=> unnötige Kontrollen vermeiden

Ziele des Meilensteins (2)

Modellentwicklung in mehreren Schritten:

- Einfache Regeln für offensichtliche Betrugsfälle
- Klassifikator für **Betrugswahrscheinlichkeit**
- Regressionsmodell zur **Schadenshöhe**
- Kombination beider Modelle in **Entscheidungslogik**
- **Threshold- & Sensitivitätsanalyse** zur Strategieoptimierung

Konkrete Handlungsempfehlungen für den operativen Einsatz

Anforderungen an Analyseverfahren

- Mehr als nur Modellgüte: Entscheidungskriterien im Praxiseinsatz
- **weitere zentrale Anforderungen** gleichrangig berücksichtigt u.a.:
 - **Verständlichkeit:** Ergebnisse nachvollziehbar & visualisierbar
 - **Umsetzbarkeit:** Einfach in der Praxis einsetzbar
 - **Reproduzierbarkeit:** Konsistente Ergebnisse mit gleichem Code/Daten
 - **Skalierbarkeit:** Einsetzbar in allen Filialen, nachtrainierbar
 - **Robustheit:** Stabil bei Datenschwankungen & erneutem Training

Merkmalsraum der Analysedaten (1)

1. Merkmalsquellen (Datenursprung):

- **Originaldaten** (z. B. Zahlungsart, Uhrzeit, Produktkategorie)
- **Automatisch generierte Systemdaten** (z. B. Kamerasystem, Rückmeldungen)
- **Berechnete Merkmale** (z. B. Transaktionsdauer, Preisabweichung)

→ Kombination ermöglicht Erkennung von Mustern als auch komplexer Zusammenhänge.

Merkmalsraum der Analysedaten (2)

2. Thematische Kategorien (inhaltliche Gruppierung):

- **Kundenverhalten** (z. B. Scanverhalten, Rückmeldungen)
 - **Preis & Rabattnutzung**
 - **Kamerabasierte Hinweise**
 - **Zeitliche Informationen**
 - **Produktbezogene Angaben**
- ermöglicht Strukturierung von Mustern für die Modellierung

Modellbildungsprozess & Datenkategorien

Kategorie	Anzahl Datensätze	NORMAL	FRAUD	Anteil FRAUD (%)	Gesamtschaden (€)
Unscanned	377	0	377	100,0 %	5.088 €
Fehlerhafte-Rabatte	1.521	0	1.521	100,0 %	11.058 €
Übrige Rabatte	9.562	8.401	1.161	12,15 %	7.960 €
Übrige	136.564	134.968	1.596	1,17 %	11.057 €
Gesamt	148.024	143.369	4.655	3,15 %	35.163 €

Modellbildungsprozess & Datenkategorien

- **Rabattsystematik auffällig**, aber schwer übertragbar auf neue Filialen → aktuell keine starren Rabattregeln implementiert.
- Mögliche technische Prävention:
 - **Rabattfunktion** bei nicht rabattfähigen Produkten **deaktivieren**
 - Nutzung vordefinierter **Rabatt-Barcodes** zur Kontrolle
 - Ergänzend: **statische Modellregeln** für auffällige Rabattmuster definieren (s. auch **Abschnitt 5**)

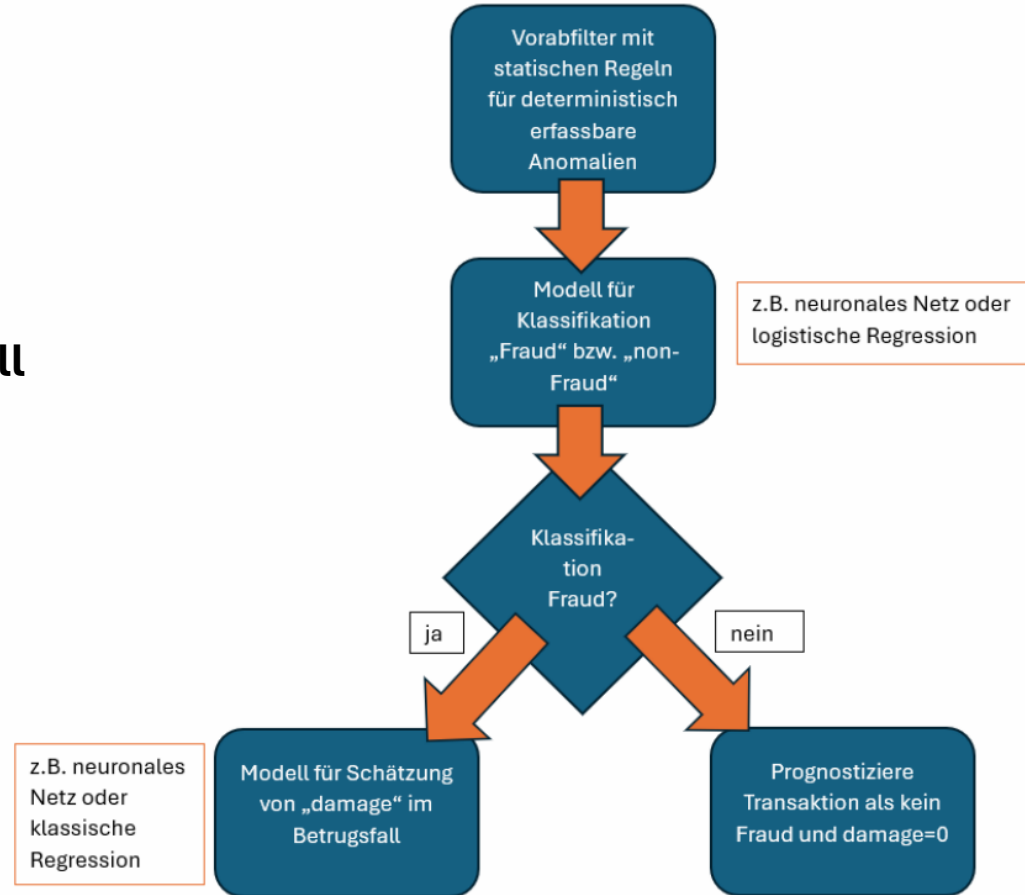
2. Statische Regeln

Erinnerung: Unser dreistufiges Modell

1. Statische Regeln

2. Klassifikationsmodell

3. Regressionsmodell



Statische Regeln zur Vorfilterung (1)

- Ziel: einfache, **interpretierbare Regeln mit hoher Präzision bei minimaler Komplexität**
- **Methodik:**
 - Daten kategorial / binär kodiert
 - Analyse von Regeln mit ein bis zwei Merkmalen, um Überanpassung zu vermeiden und Interpretierbarkeit zu gewährleisten
 - Bewertung: Güte der Vorhersage höher als bei dem anschließenden Klassifikationsmodell

Statische Regeln zur Vorfilterung (2)

- Regeln mit einem Merkmal als Basis:
 - **has_unscanned** == True mit perfekter Vorhersage von Betrugsfällen
 - **has_missing** == True ebenfalls mit perfekter Vorhersage
- Wirtschaftlicher Nutzen > **5.000 €** potenziell verhindert, aber nur **geringe Abdeckung** der gesamten Fälle (400)

	Regel	TP	FP	FN	TN	Precision	Recall	FPR	FNR	Verhinderter Schaden
0	has_unscanned == True	377	0	4278	143369	1.0	0.080988	0.0	0.919012	5088.38
1	has_missing == True	16	0	4639	143369	1.0	0.003437	0.0	0.996563	200.07

Bewertung der statischen Regeln

- **Regeln mit zwei Merkmalen** enthalten entweder wiederum has_unscanned oder haben eine FPR von über 50% und sind daher **nicht sinnvoll**; nur Verwendung der beiden Einzelregeln
- Einzelregeln sehr präzise und ideal für vorgesehenen Einsatz
- has_unscanned & has_missing: FPR = 0, d. h. kein einziger False Positive Fall im Training
- **Ggf. Erweiterung** um zuvor besprochene statische Regeln gegen **Rabattbetrug**

3. Klassifikation

Klassifikation der Transaktionen

- Klassifikationsmodell liefert **Score zwischen 0 und 1** je Transaktion (→ Fraud- „Wahrscheinlichkeit“)
- Technisch keine echten Wahrscheinlichkeiten, aber gut interpretierbare Scores (nach Kalibrierung)
- Entscheidung erfolgt über einen **Threshold** (z. B. 0.5): Ab diesem Wert wird als FRAUD klassifiziert

Modellentwicklung & Evaluation

- **Iterativer Prozess** mit gezielter Auswahl leistungsfähiger Klassifikationsmodelle
- **4 zentrale Schritte:**
 - **Modellauswahl** & Vorabtests → ungeeignete Modelle ausgeschlossen
 - **Hyperparameter-Optimierung** & Kalibrierung der Scores
 - **Merkmalsauswahl** zur Reduktion & Robustheit
 - **Evaluation** mit mathematischen Verfahren & betriebswirtschaftlicher Bewertungsfunktion

Verglichene Modellklassen

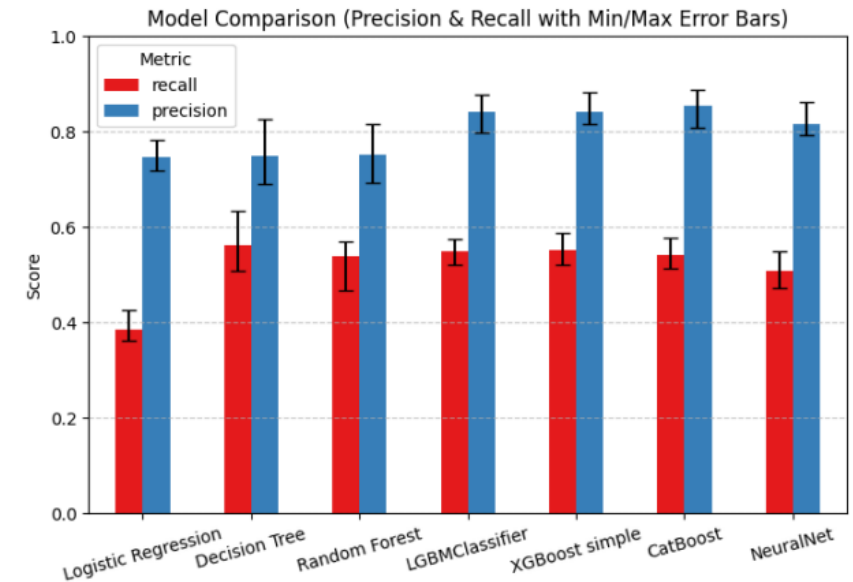
- **Baseline-Modelle:** Logistische Regression & Entscheidungsbaum
- **Fortgeschrittene Modelle:** Random Forest, Boosting (XGBoost, CatBoost), neuronale Netze
- **Boosting-Modelle** performten am besten → gezielte Weiterentwicklung
- **Neuronale Netze** zeigten gute Einzelresultate, aber instabil & sensitiv gegenüber Parametern

Modellvergleich (1)

- Trainingsdaten ohne durch statische Regeln eindeutig erkannte Fälle
- Lineares Modell: nur **3 ausgewählte Merkmale** aus Phase 2
- Alle anderen Modelle: **29 gezielt ausgewählte Features** (z. B. Zahlungsmittel, Tageszeit, Kamerasignale)
- Präprozessierung: One-Hot-Encoding + ggf. Skalierung
- Bewertung mit 5×5-facher Kreuzvalidierung unter Beibehaltung der Klassenverteilung (stratified CV)

Modellvergleich (2)

- **XGBoost, CatBoost, LightGBM** liefern beste Ergebnisse – sowohl statistisch als auch wirtschaftlich
 - **≈ 55 % der Fraud-Fälle erkannt**, hohe Präzision → wenige unnötige Kontrollen
 - +5 % Recall gegenüber neuronalen Netzen
 - **Logistische Regression deutlich schlechter** bei Recall & Schadenserkennung



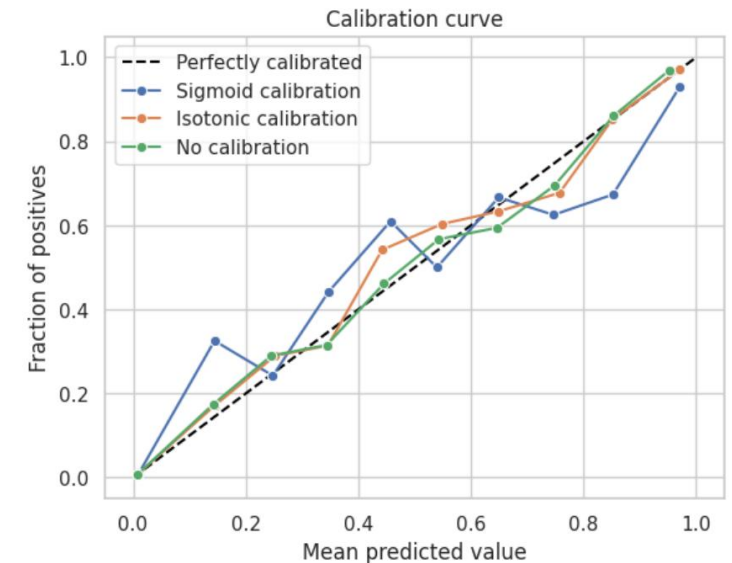
	precision	recall	f1	auc-pr	damage_prevented	Bewertung
Logistic Regression	0.746	0.385	0.508	0.431	2219.208	-3271.956
Decision Tree	0.749	0.561	0.641	0.655	3648.474	-1585.290
Random Forest	0.753	0.540	0.628	0.681	3484.541	-1748.223
LGBMClassifier	0.843	0.549	0.664	0.729	3524.914	-1020.850
XGBoost simple	0.842	0.552	0.667	0.730	3555.049	-982.715
CatBoost	0.854	0.543	0.664	0.733	3510.251	-978.913
NeuralNet	0.816	0.508	0.626	0.681	3356.139	-1468.625

Modellvergleich (3)

- **XGBoost mit besserem Recall, CatBoost mit höherer Präzision** → Trade-off zwischen Entdeckungsrate und Kontrollkosten
- Random Forest unterliegt dem Einzelbaum – trotz Theorievorteil
 - Ursache: fehlende Hyperparameter-Optimierung
- CatBoost leicht besser, aber **Entscheidung zugunsten XGBoost** aus praktischen Gründen:
 - Starke Verbreitung
 - Gute Dokumentation
 - Effizientes Training
 - Besser wart- & erweiterbar im operativen Einsatz
- **Erfüllt alle Anforderungen:** Verständlichkeit, Skalierbarkeit, Robustheit, Reproduzierbarkeit

Kalibrierung & Schwellenwertwahl

- XGBoost-Scores \neq echte Wahrscheinlichkeiten, von daher sollte im Nachgang der Entscheidungsschwellwert (Sicherheit des Modells, dass Betrug vorliegt) kalibriert werden
- Aber: **Modell zeigt ohne Kalibrierung gute Performance.** Performance sogar besser als bei nachträglicher Schwellwertoptimierung!
- Standardwert 0.5 liefert stabilere & bessere Ergebnisse \rightarrow keine Schwellenanpassung nötig



4. Regression

Schadensschätzung per Regression

- Ziel: **Schadenshöhe im Betrugsfall** prognostizieren – unabhängig von der Klassifikationssicherheit
- Ergänzt die Klassifikation um quantitative Risikoabschätzung pro Transaktion
- Ermöglicht differenzierte Entscheidungen: z. B. Kontrolle trotz niedriger Score-Wahrscheinlichkeit bei hohem vermutetem Schaden
- **Boosting-Modelle** erneut am besten, verwendet wird ein XGBoost-Regressor

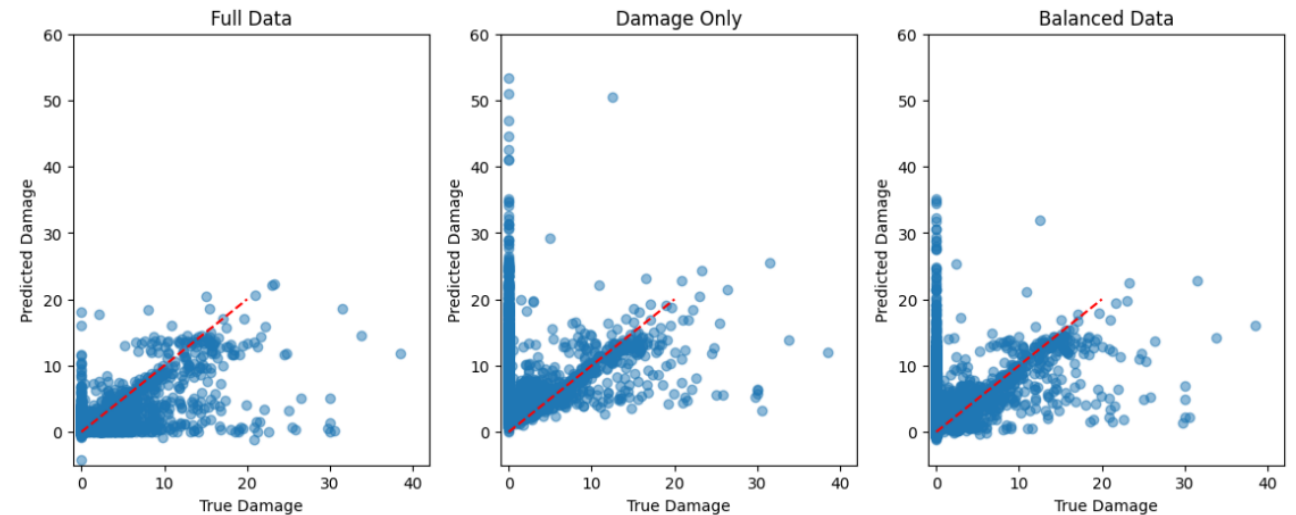
Trainingsvarianten (1)

- 3 Trainingsvarianten zur Modellierung getestet:
 - **Vollständiger Datensatz** (inkl. Schaden = 0): realistisch, aber stark unausgewogen
 - **Balanced Set** (gleich viele Schaden / kein Schaden): sensitiv, aber nicht repräsentativ
 - **Nur Schadensfälle**: genauer für Betrug, aber nicht einsetzbar bei normalen Transaktionen
- Zielkonflikt: Generalisierung vs. Präzision vs. Repräsentativität

Trainingsvarianten (2)

- Alle Varianten haben Schwierigkeiten bei der Vorhersage hoher Schäden
- Im Bereich 0–10 €: hohe Streuung, quadratische Verteilung
- „Nur-FRAUD“-Modell überschätzt normale Transaktionen stark, trifft aber hohe Schäden am besten
- Klassische Metriken (R^2 , MAE etc.) nur bedingt aussagekräftig im Vergleich
- „Damage-only“-Variante versagt bei Generalisierung, balanced liegt dazwischen

Resultate



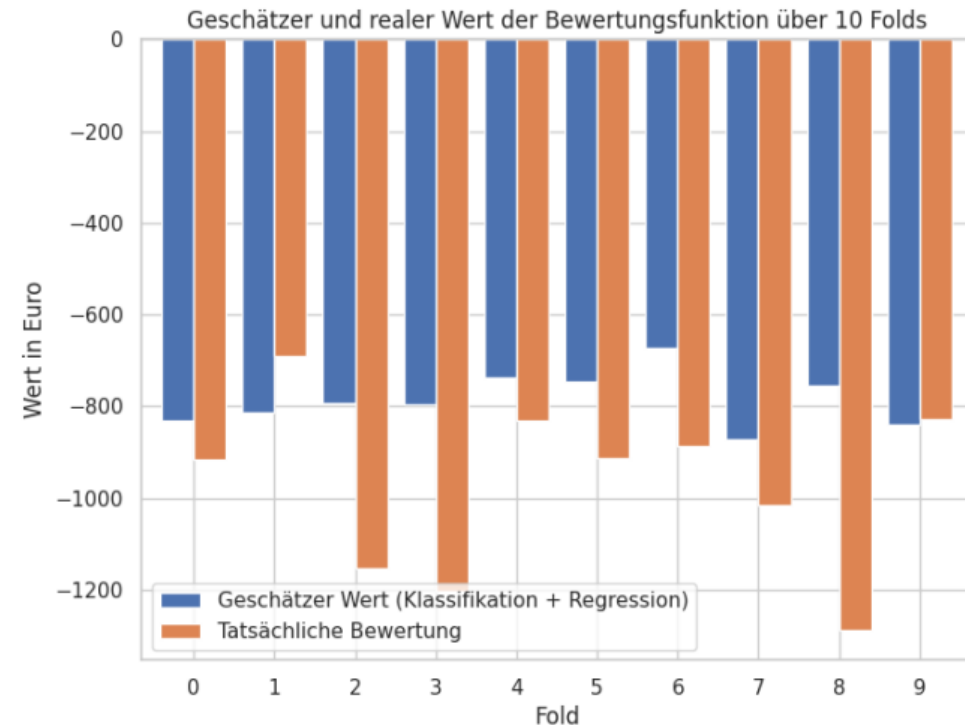
5. Gesamtmodell

Simulierte Bewertungsfunktion

- Ziel: wirtschaftlich sinnvolle Kontrollentscheidung für jede Transaktion
- Vergleich zweier Szenarien:
 - **Keine Kontrolle** → potenzieller Schaden bei nicht erkanntem Betrugsfall: $P(\text{FRAUD}) * \text{erwarteter Schaden}$
 - **Kontrolle** → Mischung aus erwarteter Fraud-Prämie (bei richtiger Klassifikation) & False-Positive-Kosten (bei Falschklassifikation): $P(\text{FRAUD}) * 5 \text{ €} - P(\text{NORMAL}) * 10 \text{ €}$
- **Kombiniertes Modell (Klassifikation + Regression)** simuliert Entscheidung für gesamten Datensatz
 - Wahrscheinlichkeiten $P(\text{FRAUD})$ bzw. $P(\text{NORMAL})$ aus Klassifikationsmodell
 - Erwartungswert des Schadens $E(\text{Schaden})$ aus Regressionsmodell

Simulierte Bewertungsfunktion

- Nur das Modell auf **vollständigem Datensatz** kann die tatsächliche Bewertungsfunktion realitätsnah approximieren
- Andere Varianten (balanced / damage-only) liefern massiv verzerrte Werte (falsche Mittelwerte: 5.1 / 4.7 statt 0.21)
- **Modell ist leicht optimistisch**, aber klar robustester Ansatz für praxisnahe Entscheidungen



Zusätzliche Optionen im Modell

- Neben den aktuell fixen Werten als **Belohnung bzw. Bestrafung** für richtig erkannte bzw. fälschlich als Betrug markierte Transaktionen können die Werte anhand der **Konfigurationsdatei** „beliebig“ **verändert** werden
- Die für **Rabatte ausgeschlossenen Produktkategorien** könnten ebenfalls per Konfiguration angepasst werden, sodass Rabatte in Kombination mit diesen Produktkategorien als Betrug (statisch) bewertet werden

Schaden durch Rabattbetrug

- Modell erkennt **62 % der unberechtigten Rabattfälle**
- Dadurch können im Schnitt 63 % des zugehörigen Schadens verhindert werden
- → **Hohes Präventionspotenzial** bei Rabattmissbrauch

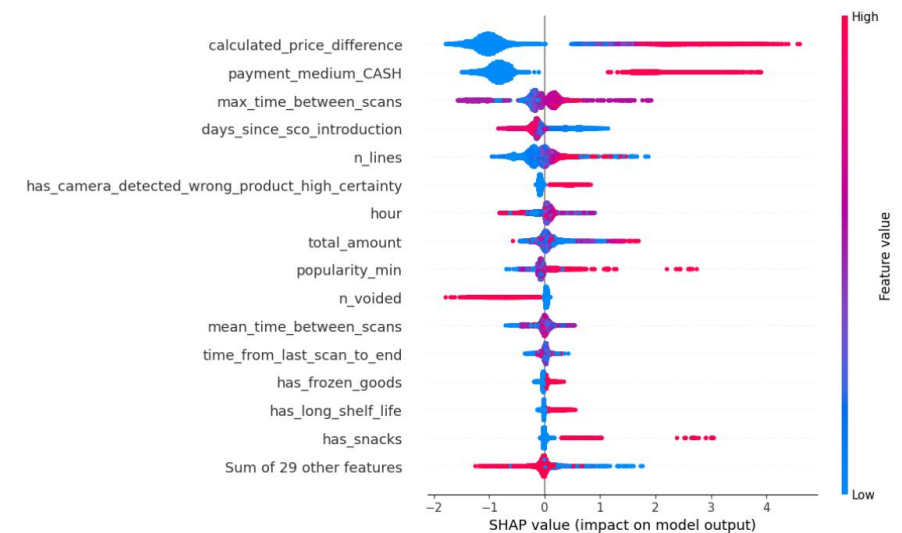
6. Bewertung & Empfehlung

Wirtschaftlicher Mehrwert des Modells

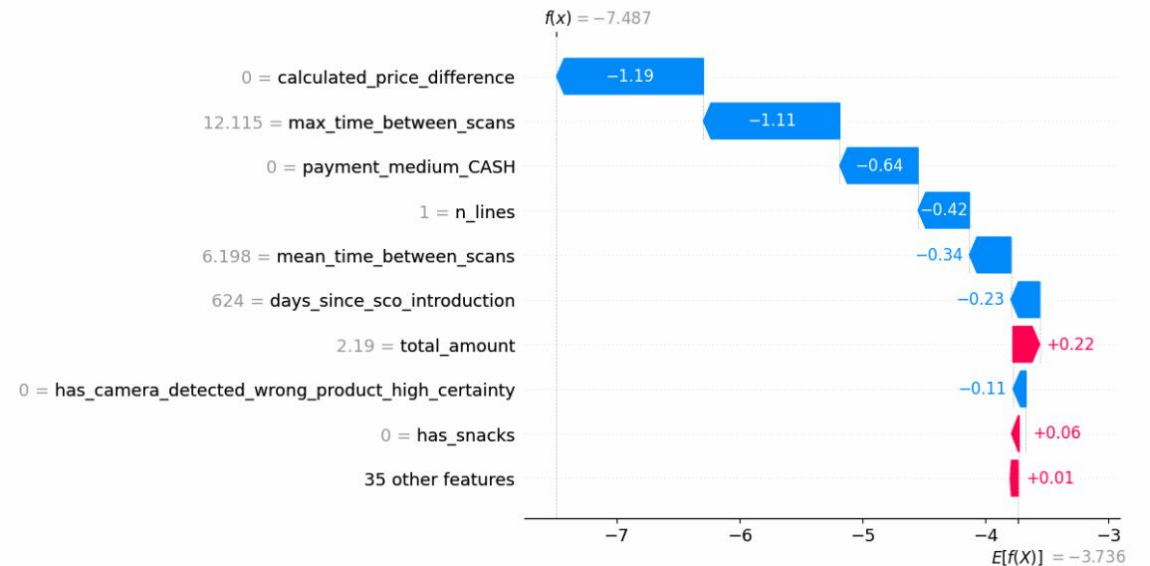
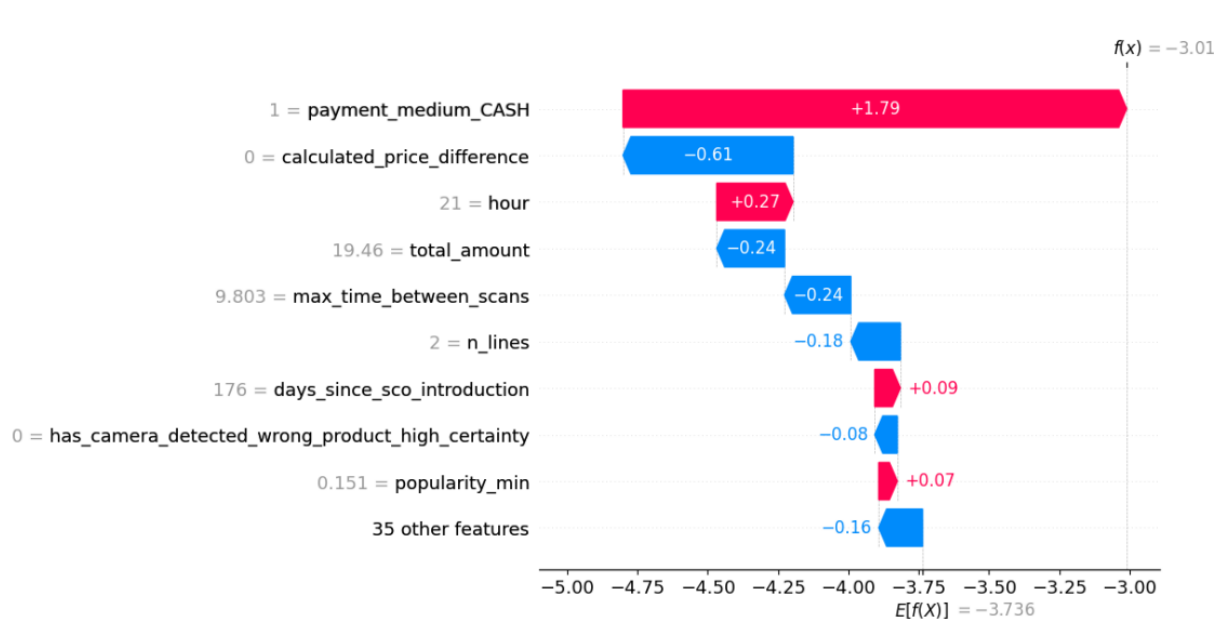
- Umfangreiche Evaluation: 200 Testläufe (5×CV mit 40 Wiederholungen)
- **Durchschnittlicher Mehrwert: 0,22 € pro Transaktion** (nach Bewertungsfunktion)
- Betrachtet man nur den verhinderten Schaden, also ohne Abzug der Kontrollkosten und ohne Bonus für entdeckte FRAUD-Fälle, ergibt sich ein mittlerer Wert von **0,15 € pro Transaktion**.
- Ergebnis gilt als robuste, belastbare Schätzung der Modellleistung
- Bezieht sich auf Testdaten (20 % des Gesamtbestands)

Sensitivitätsanalyse: Einflussfaktoren im Modell (1)

- Wichtigste Prädiktoren:
 - **Bargeldzahlung** = stärkster Einzelindikator für Fraud
 - **Preisabweichung & Kamera-Hinweise** erhöhen Risiko deutlich
- Zeitliche Merkmale (Scan-Dauer, Tageszeit) mit moderatem Einfluss
- **Jüngere Kamerasysteme** liefern **weniger aussagekräftige Daten**
- Modell trifft **nachvollziehbare Entscheidungen**, keine Blackbox

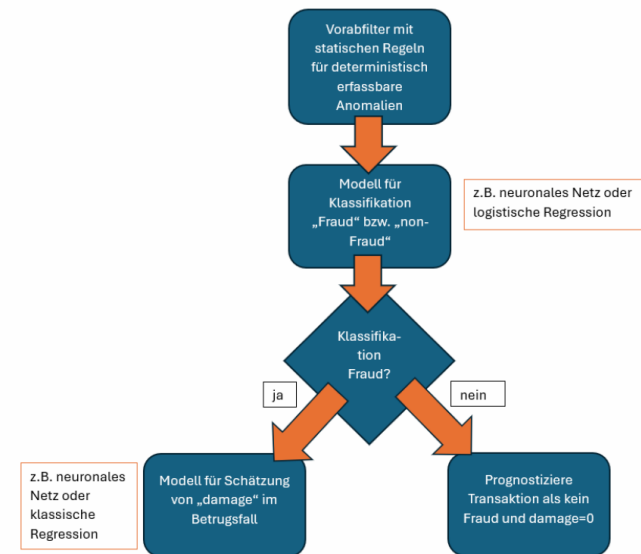


Sensitivitätsanalyse: Einflussfaktoren im Modell (2)



Handlungsempfehlungen & Modellpflege

- **Gesamtmodell** (Statische Regeln + Klassifikation + Regression) ist betriebsreif → sollte **diskretionäre** Kontrollen durch **datenbasierte Entscheidungen** ersetzen
- Ergebnis: Schäden vermeiden & Personal effizienter einsetzen
- Regelmäßige **Rekalibrierung** empfohlen bei:
 - neuen Filialen / Sortimenten
 - veränderten Kundengruppen
 - neuen Kameradaten / Kontrollrückmeldungen



7. Ausblick

Technische Umsetzung

- Modell wird über eine **REST-API** ins Kassensystem integriert
 - Input: Transaktionsdaten (JSON)
 - Output: Echtzeitentscheidung + Schadenprognose + Begründung
- **Codeversionierung & Nachtraining** über GitHub möglich
- **Evaluation** mit Echtdateien der Wertkauf GmbH geplant
- Langfristige **Erweiterung denkbar**: z. B. durch Kundenhistorie, Treuekarten, Warenkorbinhalte



**Vielen Dank für Ihre
Aufmerksamkeit!**

Fragen & Anregungen?