

Detekcja oszustw z wykorzystaniem metod wrażliwych na koszt

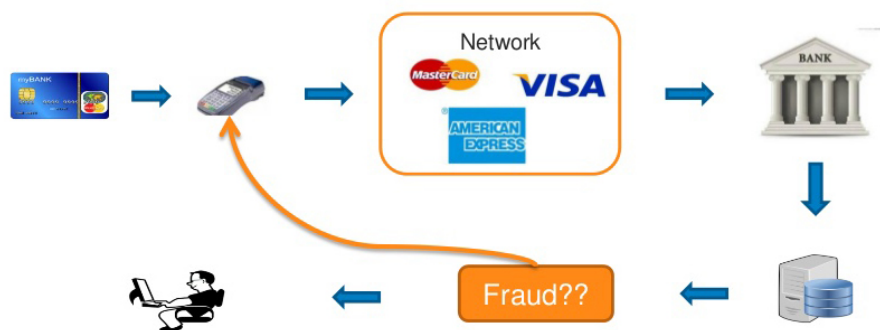
Patryk Wielopolski

24 grudnia 2019

Rozdział 1

Wstęp

Wraz z postępującym rozwojem technologii w ostatnich latach, który umożliwia Nam posługiwanie się kartami kredytowymi w codziennych transakcjach, podczas zakupów w internecie, a także ostatnimi czasy również i za pomocą urządzeń elektronicznych takich jak telefony komórkowe lub inteligentne zegarki, znacząco wzrosła ilość transakcji dokonywanych elektronicznie. Niewątpliwie zwiększyło to wygodę dokonywania zakupów, natomiast z drugiej strony spowodowało to istotny wzrost możliwości dokonywania przestępstw wynikających z tego rodzaju płatnościami, tj. przekazania numeru karty nieznajomym, utraty lub kradzieży karty, kradzieży przesyłki lub skopiowania danych karty, które zwykle kończą się pojawieniem się na koncie nieautoryzowanych płatności.



Rysunek 1.1: Schemat przepływu informacji związanych z płatnością kartą kredytową.

Banki oraz instytucje finansowe w celu przeciwdziałania takim działaniom

umieszczają specjalistyczne systemy detekcji oszustw w procesie dokonywania płatności, który jest przedstawiony schematycznie na Rysunku 1. W momencie przyłożenia karty kredytowej do terminala następuje połączenie z operatorem karty, który następnie łączy się z bankiem klienta. Kolejno w odpowiednich bazach danych sprawdzane są dostępne środki, a także poprawność transakcji. W tym momencie działa nasz system, który w razie wykrycia niezgodności odrzuca transakcję, a sprawa jest kierowana do analityka, który bliżej przygląda się sprawie.

Jak możemy zauważyć system do detekcji oszustw stanowi jedynie pewien element całego schematu. Do niedawna najbardziej popularnym sposobem do wykrywania oszustw były systemy oparte na regułach biznesowych. Przykładowe reguły, które mogą być używane w takich modelach mogą przyjmować następującą postać:

- Więcej niż 4 wypłaty z bankomatu w ciągu godziny;
- Więcej niż 2 transakcje w ciągu 5 minut;
- Transakcja w sklepie karta a następna zaraz w internecie.

W przypadku, gdy rozpatrywana transakcja spełnia chociaż jedną z reguł to jest ona blokowana i użytkownik jest informowany o odrzuceniu transakcji. To podejście jest stosunkowo łatwe do implementacji oraz jest bardzo proste do interpretacji, lecz niestety oszuści wraz z upływem czasu wpadają na nowe pomysły w jaki sposób można dokonywać oszustw, natomiast systemy złożone ze sztywnych reguł nie są w stanie wykrywać nowych zachowań i generować coraz to nowszych reguł. Tą niedogodność są w stanie rozwiązać model uczenia maszynowego, które tworzą pewnego rodzaju reguły, które nie są bezpośrednio wprowadzane przez użytkownika, lecz automatycznie wykrywane przez algorytmy na podstawie dostarczonych danych historycznych. Wraz z dynamicznym rozwojem tej dziedziny nauki zaczęły powstawać coraz bardziej wyrafinowane modele, które ostatecznie są wykorzystywane w produkcyjnych środowiskach instytucji finansowych.

Rozdział 2

Wprowadzenie teoretyczne

W tej części zostaną wprowadzone wszelkie potrzebne miary skuteczności modeli oraz modele predykcyjne, które zostaną wykorzystane do przeprowadzenia eksperymentu.

Modele predykcyjne zostały podzielone na dwie kategorie: standardowe oraz wrażliwe na koszt. Pierwsze z nich są powszechnie wykorzystywane w standardowych aplikacjach. Drugie z nich dzielą się na dwie podkategorie - *Cost Sensitive Training* oraz *Cost Sensitive Classification*.

2.1 Miary skuteczności modeli

2.1.1 Macierz pomyłek

Jedną z najpopularniejszych metod wykorzystywanych do oceny skuteczności modeli w zadaniach klasyfikacji binarnej jest macierz pomyłek. Oznaczone dane zestawiamy z klasami nadanymi przez model, a następnie sprawdzamy czy poprawnie zaklasyfikowaliśmy poszczególne obserwacje. W ten sposób powstają cztery możliwe sytuacje, które są przedstawione w Tabeli 2.1:

- Prawdziwie pozytywna klasyfikacja (*ang. True Positive (TP)*) - obserwacja była pozytywna i zaklasyfikowaliśmy ją jako pozytywną
- Fałszywie negatywna klasyfikacja (*ang. False Negative (FN)*) - obserwacja była pozytywna i zaklasyfikowaliśmy ją jako negatywną. W nomenklaturze statystycznej nazywana również błędem drugiego rodzaju.
- Fałszywie pozytywna klasyfikacja (*ang. False Positive (FP)*) - obserwacja była negatywna i zaklasyfikowaliśmy ją jako pozytywną. W nomenklaturze statystycznej nazywana również błędem pierwszego rodzaju.
- Fałszywie negatywna klasyfikacja (*ang. False Positive (FN)*) - obserwacja była pozytywna i zaklasyfikowaliśmy ją jako negatywną

		Predykcja	
		Pozytywna	Negatywna
Prawda	Pozytywna	TP	FN
	Negatywna	FP	TN

Tabela 2.1: Macierz pomyłek

Dodatkowo na podstawie macierzy pomyłek możemy zdefiniować następujące miary skuteczności modeli:

Skuteczność, która określa jaki procent wszystkich obserwacji zaklasyfikowaliśmy poprawnie.

$$\text{Skuteczność} = \frac{TP + TN}{TP + FP + FN + TN}$$

Precyzja, która mierzy ile procent zaklasyfikowanych próbek jako pozytywne faktycznie były pozytywne.

$$\text{Precyzja} = \frac{TP}{TP + FP}$$

Czułość, która mówi Nam o procencie znalezionych wszystkich pozytywnych przypadków.

$$\text{Czułość} = \frac{TP}{TP + FN}$$

F1 Score, który jest średnią harmoniczną precyzji oraz czułości.

$$\text{F1 Score} = 2 \cdot \frac{\text{Precyzja} \cdot \text{Czułość}}{\text{Precyzja} + \text{Czułość}}$$

Wszystkie wyżej wspomniane miary skuteczności przyjmują wartości z zakresu $[0, 1]$ oraz wyższa wartość oznacza lepszy model.

2.1.2 Miary skuteczności modeli wrażliwe na koszt

Motywacją do powstania miar wrażliwych na koszt jest rzeczywista ewaluacja modeli. Podstawowe metryki nie uwzględniają różnicy w kosztach pomyłki dla FP i FN. Ponadto koszt fraudów znacząco różni się w zależności od przypadku.

W celu wprowadzenia potrzebnych miar skuteczności potrzebujemy wprowadzić tzw. macierz kosztu, która jest zaprezentowana w Tabeli 2.2, gdzie poszczególne komórki oznacza odpowiadające wartość kosztu predykcji. Następnie definiujemy następującą wartość:

$$\text{Koszt}(f(\mathbf{x}_i^*)) = y_i(c_i C_{TP_i} + (1 - c_i) C_{FN_i}) + (1 - y_i)(c_i C_{FP_i} + (1 - c_i) C_{TN_i}),$$

gdzie

		Predykcja	
		Oszustwo	Normalna
Prawda	Oszustwo	C_{TP_i}	C_{FN_i}
	Normalna	C_{FP_i}	C_{TN_i}

Tabela 2.2: Macierz kosztu

- $\mathbf{x}_i^* = [\mathbf{x}_i, C_{TP_i}, C_{FP_i}, C_{FN_i}, C_{TN_i}]$ - wektor zawierający wartości cech i-tej obserwacji rozszerzony o koszt klasyfikacji
- C - koszt klasyfikacji i-tej obserwacji
- $f(\cdot)$ - model predykcyjny
- y_i - prawdziwa wartość i-tej obserwacji
- c_i - predykcja dla i-tej obserwacji

Następnie wprowadzamy następujące miary skuteczności modeli:

$$\text{Koszt całkowity}(f(\mathbf{S})) = \sum_{i=1}^N \text{Cost}(f(\mathbf{x}_i^*)) \quad (2.1)$$

$$\text{Oszczędności} = \frac{\text{Koszt}_l(\mathbf{S}) - \text{Koszt}(f(\mathbf{S}))}{\text{Koszt}_l(\mathbf{S})} \quad (2.2)$$

gdzie

- \mathbf{S} - data set
- $\text{Koszt}_0 = \min\{\text{Cost}(f_0(\mathbf{S})), \text{Koszt}(f_1(\mathbf{S}))\}$
- $f_a(\mathbf{S}) = \mathbf{a}$ gdzie $a \in \{0, 1\}$

Wartość $f_a(\mathbf{S})$ możemy rozumieć jako przypadek naiwnego modelu, który wszystkim obserwacjom przyznaje wartość a . Natomiast Koszt_0 oznacza wybranie naiwnego klasyfikatora, który generuje mniejsze koszty. Zatem ostatecznie Oszczędności możemy rozumieć jako procentową wartość o ile testowany model jest lepszy od naiwnego klasyfikatora.

2.2 Standardowe modele

2.2.1 Regresja logistyczna

Regresja logistyczna należy do jednego z najbardziej podstawowych modeli statystycznych używanych do problemów klasyfikacyjnych. Jest to model, w którym

zakładamy, że logarytm szans jest kombinacją liniową predyktorów, co możemy zapisać jako:

$$\log \frac{p}{1-p} = \sum_{j=1}^k \theta^{(j)} x_i^{(j)}$$

Po nałożeniu odpowiednich przekształceń ostatecznie otrzymujemy, że estymowane prawdopodobieństwo możemy zapisać w następującej postaci

$$\hat{p} = P(y = 1 | \mathbf{x}_i) = h_{\theta}(\mathbf{x}_i) = g\left(\sum_{j=1}^k \theta^{(j)} x_i^{(j)}\right) \quad (2.3)$$

gdzie funkcja $g(z)$ to sigmoidalna funkcja łącząca, która jest opisana wzorem:

$$g(z) = \frac{1}{(1 + e^{-z})}.$$

W celu znalezienia odpowiednich parametrów definiujemy standardowa funkcja straty, która przyjmuje następującą postać:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N J_i(\theta),$$

gdzie $J_i(\theta)$ często w literaturze nazywana jest entropią krzyżową i przyjmuje postać:

$$J_i(\theta) = -y_i \log(h_{\theta}(\mathbf{x}_i)) - (1 - y_i) \log(1 - h_{\theta}(\mathbf{x}_i)) \quad (2.4)$$

Następnie wykorzystując algorytm optymalizacyjny, który będzie w stanie znaleźć minimum funkcji straty, znajdujemy odpowiednie parametry modelu.

Korzystając z Rysunku 2.2.1 możemy zauważyć, że koszt klasyfikacji przyjmuje następujące wartości:

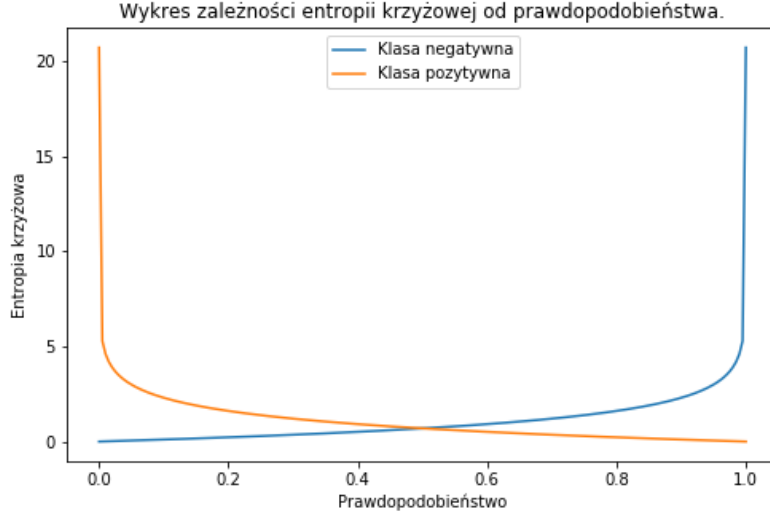
$$J_i(\theta) \approx \begin{cases} 0, & \text{if } y_i \approx h_{\theta}(\mathbf{x}_i), \\ \infty, & \text{if } y_i \approx (1 - h_{\theta}(\mathbf{x}_i)). \end{cases}$$

To znaczy, że w przypadku prawidłowego zaklasyfikowania danej obserwacji funkcja kosztu przyjmuje wartość bliską zeru, natomiast w przypadku pomyłki nieskończoności. Zatem stąd możemy wywnioskować, że koszty popełnienia bądź niepopelnienia błędu są następujące:

$$C_{TP_i} = C_{TN_i} \approx 0$$

$$C_{FP_i} = C_{FN_i} \approx \infty$$

Porównując otrzymane rezultaty z macierzą kosztów możemy zauważyć, że te wartości znacząco odbiegają od tego co w rzeczywistości chcemy otrzymać. Stąd powstaje motywacja, aby zmodyfikować podaną funkcję celu i stworzyć regresję logistyczną wrażliwą na koszt, która będzie optymalizować właściwe dla Nas wartości.



Rysunek 2.1: Wykres zależności funkcji entropii krzyżowej od prawdopodobieństwa predykcji danej obserwacji. Kolor niebieski przedstawia wykres dla próbki o prawdziwej klasie pozytywnej, natomiast kolor pomarańczowy dla negatywnej.

2.2.2 Drzewo decyzyjne

Drzewo klasyfikacyjne to przykład jednego z rodzaju drzew decyzyjnych, którego celem jest znalezienie najlepszego rozróżnienia pomiędzy klasami. W ogólności drzewo decyzyjne składa się z zestawu reguł.

Drzewo składa się z węzłów, które są reprezentowane przez parę (\mathbf{x}^j, l^j) , która oznacza podział zbioru obserwacji \mathcal{S} na dwa zbiory: \mathcal{S}^l oraz \mathcal{S}^r względem wektora obserwacji \mathbf{x} oraz progu decyzyjnego l^j w następujący sposób:

$$\mathcal{S}^l = \{\mathbf{x}^* : \mathbf{x}^* \in \mathcal{S} \wedge x_i^j \leq l^j\},$$

$$\mathcal{S}^r = \{\mathbf{x}^* : \mathbf{x}^* \in \mathcal{S} \wedge x_i^j > l^j\},$$

gdzie x_i^j reprezentuje i -ty atrybut wektora \mathbf{x} . Ponadto l^j jest wartością taką, że $\min \mathbf{x}^j \leq l^j < \max \mathbf{x}^j$. Ponadto warto zauważyć, że $\mathcal{S}^l \cup \mathcal{S}^r = \mathcal{S}$, co oznacza, że nasz podział rozdziela wektor obserwacji na dokładnie dwa rozłączne zbiory. Po znalezieniu optymalnego podziału zliczamy ilość pozytywnych próbek:

$$\mathcal{S}_1 = \{\mathbf{x}^* : \mathbf{x}^* \in \mathcal{S} \wedge y_i = 1\},$$

a następnie zliczamy procent pozytywnych próbek jako:

$$\pi_1 = \frac{|\mathcal{S}_1|}{|\mathcal{S}|}.$$

Następnie dla każdego z liści jest obliczana wielkość jego zanieczyszczenia.

- Misclassification: $I_m(\pi_1) = 1 - \max(\pi_1, 1 - \pi_1)$
- Entropy: $I_e(\pi_1) = -\pi_1 \log(\pi_1) - (1 - \pi_1) \log(1 - \pi_1)$
- Gini: $I_g(\pi_1) = 2\pi_1(1 - \pi_1)$

Następnie dla każdego proponowanego podziału dla danej reguły (\mathbf{x}^j, l^j) liczony jest przyrost czystości w następujący sposób:

$$\text{Gain}(\mathbf{x}^j, l^j) = I(\pi_1) - \frac{|\mathcal{S}^l|}{|\mathcal{S}|} I(\pi_1^l) - \frac{|\mathcal{S}^r|}{|\mathcal{S}|} I(\pi_1^r),$$

gdzie $I(\pi_1)$ może być dowolną z zaproponowanych miar zanieczyszczenia. Ostatecznie wybiera się ten podział, który maksymalizuje przyrost czystości, a następnie dzieli się zbiór \mathcal{S} na podzbiory \mathcal{S}^l i \mathcal{S}^r . W taki sposób jest pokonywany pojedynczy podział zbioru dla konkretnego węzła. Całe drzewo tworzy się poprzez kolejne podziały węzłów aż do momentu dotarcia przez algorytm do kryterium stopu.

2.2.3 Las losowy

Kolejne modele są przedstawicielami szerokiej klasy metod *ensemble*, których celem jest złożenie predykcji wielu klasyfikatorów bazowych, aby poprawić generalizację pojedynczego estymatora. Wśród nich wyróżniamy dwie główne kategorie. Pierwsza z nich to metody uśredniania, które polegają na zbudowaniu wielu niezależnych klasyfikatorów, a następnie uśrednianie wyników predykcji. Przedstawicielem tej kategorii jest las losowy. Natomiast druga polega na iteracyjnym budowaniu kolejnych modeli, które próbują zredukować obciążenie poprzednika. Bardzo powszechnie znanym reprezentantem jest algorytm XGBoost, którym zajmujemy się w następnym podrozdziale.

Metody *ensemble* wykorzystują metody próbkowania w celu utworzenia różnych klasyfikatorów bazowych, aby następnie dokonać ostatecznej predykcji. Metody te są używane, aby zredukować wariancję klasyfikatora bazowego (zazwyczaj drzewa decyzyjnego) poprzez losowe dobieranie próbki i/lub zmiennych, na których model będzie uczony. W wielu przypadkach stworzenie modelu opartego o *bagging* jest znacznie prostsze, ponieważ wymaga jedynie zmiany próbkowania danych bez naruszania modelu bazowego, natomiast metody typu *boosting* wymagają zmiany całego algorytmu. Z licznych obserwacji wynika, że pierwsza z metod dużo lepiej radzi sobie mają jako bazowe klasyfikatory złożone modele, w przeciwieństwie do drugiej, która zazwyczaj najlepiej performuje wykorzystując proste modele (np. płytkie drzewa decyzyjne, tzn. o małej głębokości).

Przykładowe metody losowania próbek do modeli bazowych:

- *Pasting* - losowanie obserwacji bez powtórzeń
- *Bagging* - losowanie obserwacji z powtórzeniami

- *Random Subspaces* - losowanie podzbioru zmiennych
- *Random Patches* - losowanie podzbioru zmiennych oraz obserwacji

W przypadku lasu losowego proces losowania próbek jest podzielony na dwie fazy. Pierwsza z nich polega na próbkowaniu z powtórzeniami obserwacji ze zbioru treningowego dla każdego z osobnych estymatorów bazowych. Następnie podczas fazy tworzenia kolejnych węzłów w drzewach wybierany jest losowy podzbiór zmiennych, które mogą być w tym kroku wykorzystane.

Celem tych dwóch różnych źródeł losowości jest redukcja wariancji lasu. Pojedyncze drzewa mają tendencję do zbytniego dopasowywania się do danych, zatem losowanie poszczególnych zmiennych w każdym kolejnym węźle pomaga to zredukować. Natomiast losowanie różnych próbek do każdego z klasyfikatorów pozwala na stworzenie lekko odmiennych modeli bazowych.

2.2.4 XGBoost

Tak jak wspomnieliśmy w poprzednim rozdziale algorytm XGBoost jest przykładem klasyfikatora, który w iteracyjny sposób tworzy kolejne bazowe klasyfikatory. W przypadku tego algorytmu jako klasyfikator bazowy wykorzystujemy implementację drzew decyzyjnych typu CART, które nieznacznie różnią się od standardowych drzew decyzyjnych opisywanych w 2.2.2, ponieważ liść drzewa jest rozszerzony o wartość rzeczywistą, która reprezentuje decyzję modelu. Ponieważ jest to złożenie wielu klasyfikatorów, to możemy zapisać nasz model w następującej postaci:

$$\hat{y}_i = \sum_{k=1}^K f_k(x_i), f_k \in \mathcal{F},$$

gdzie K oznacza liczbę drzew, f funkcje z przestrzeni \mathcal{F} wszystkich możliwych drzew CART. Zatem funkcją, którą będziemy optymalizować przyjmuje następującą postać:

$$\text{Obj}(\theta) = \sum_{i=1}^n l(y_i, y_i^{(t)}) + \sum_{k=1}^K \Omega(f_k) \quad (2.5)$$

W przypadku klasyfikacji przyjmujemy tak samo jak poprzednio funkcję entropii krzyżowej (2.4). Patrząc na postać modelu nie różni się ona niczym od lasu losowego. Zatem różnica między tymi modelami polega na sposobie trenowania drzew, który pokrótce opiszemy.

Ponieważ naszym modelem bazowym jest drzewo, to nie jesteśmy w stanie wprost rozwiązać zagadnienia optymalizacyjnego poprzez obliczenie gradientu funkcji i iteracyjne znalezienie rozwiązania. W tym przypadku posłużymy się treningiem addytywnym, który polega na iteracyjnym poprawianiu błędów z poprzednich modeli poprzez odpowiednie przydzielanie wag próbkom w kolejnych krokach algorytmu. Oznaczmy wartość predykcji w kroku t jako $y_i^{(t)}$.

$$y_i^{(0)} = 0$$

$$\begin{aligned}
y_i^{(1)} &= f_1(x_i) = y_i^{(0)} + f_1(x_i) \\
y_i^{(2)} &= f_1(x_i) + f_2(x_i) = y_i^{(1)} + f_2(x_i) \\
&\dots \\
y_i^{(t)} &= \sum_{k=1}^t f_k(x_i) = y_i^{(t-1)} + f_t(x_i)
\end{aligned}$$

Pozostaje zagadnienie jakie drzewo chcemy wybrać w każdym z kroków. Oczywiście takie, które optymalizuje naszą funkcję Obj. Korzystając z powyższych wzorów oraz 2.5 otrzymujemy następującą postać tej funkcji:

$$\sum_{i=1}^n l(y_i, y_i^{(t-1)} + f_t(x_i)) + \Omega(f_t) + \text{constant}.$$

Korzystając z rozwinięcia szeregu Taylora dla funkcji straty otrzymujemy ogólny wzór:

$$\text{Obj}^{(t)} = \sum_{i=1}^n [l(y_i, y_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t) + \text{constant},$$

gdzie

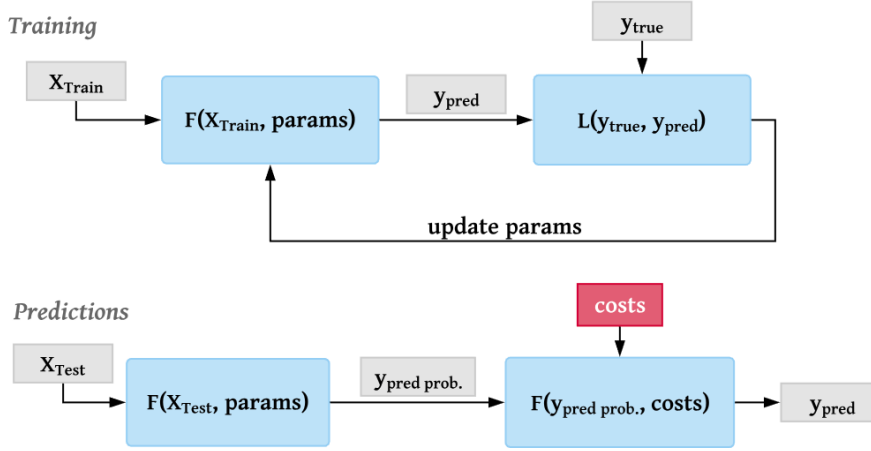
$$\begin{aligned}
g_i &= \partial_{y_i^{(t-1)}} l(y_i, y_i^{(t-1)}) \\
h_i &= \partial_{y_i^{(t-1)}}^2 l(y_i, y_i^{(t-1)})
\end{aligned}$$

Zatem po redukcji stałych, które są nieistotne z punktu widzenia optymalizacji otrzymujemy:

$$\sum_{i=1}^n [g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t)$$

2.3 Cost Dependent Classification

Metody *Cost Dependent Classification* są przykładem pierwszego rodzaju modeli wrażliwych na koszt. W przypadku tych metod trening odbywa się dopiero po etapie stworzenia modelu zwracającego prawdopodobieństwa i informacja o koszcie jest uwzględniana dopiero w tej fazie modelowania. Cały proces jest przedstawiony na Rysunku 2.2. Górna część diagramu przedstawia standardowy proces uczenia modelu predykcyjnego, natomiast dolna część reprezentuje schemat dokonywania predykcji, w którym najpierw estymujemy prawdopodobieństwa dla zbioru testowego, a następnie wykorzystując te wartości oraz koszt klasyfikacji dokonujemy ostatecznej decyzji, czy dana obserwacja jest pozytywna czy negatywna. Warto wspomnieć, że w celu dokonania kolejnego treningu modelu potrzebujemy podzbioru danych, który nie był wykorzystywany do uczenia podstawowego modelu. Najczęściej taki zbiór nazywamy walidacyjnym bądź developerskim.



Rysunek 2.2: Schemat przedstawiający proces uczenia klasyfikatora wrażliwego na koszt. Źródło: <https://towardsdatascience.com/fraud-detection-with-cost-sensitive-machine-learning-24b8760d35d9>

2.3.1 Optymalizacja progu

Metoda optymalizacji progu jest popularną metodą wyznaczania odpowiedniego progu prawdopodobieństwa powyżej którego wszystkie obserwacje oznaczamy jako pozytywnie zaklasyfikowane. Może być ona wykorzystywana nie tylko do problemów wrażliwych na koszt, lecz do dowolnie zdefiniowanego zagadnienia, w której wyznaczenie progu jest potrzebne. Jej sformułowanie wygląda następująco

$$\arg \min_{th \in [0,1]} f(y_{true}, y_{decision}),$$

gdzie

- $f(\cdot)$ - funkcja scorująca model, np. skuteczność,
- $y_{true} = (c_i)_{i=1}^n$ - wektor prawdziwych oznaczeń ,
- $y_{decision} = (p_i > th)_{i=1}^n$ - wektor binarnych odpowiedzi modelu,
- p_i - przewidywana wartość prawdopodobieństwa dla i-tej obserwacji,
- th - wartość progu decyzyjnego.

W naszym przypadku funkcja f to funkcja kosztu całkowitego (2.1). Innymi słowy, w tym przypadku będziemy szukać takiego progu decyzyjnego, który zminimalizuje sumaryczną wartość kosztów.

2.3.2 Bayesian Minimum Risk

Ryzyko poszczególnych klasyfikacji definiujemy w następujący sposób:

$$R(p_1|x_i) = L(p_1|y_1)P(p_1|x_i) + L(p_1|y_0)P(y_0|x_i),$$

$$R(p_0|x_i) = L(p_0|y_0)P(p_0|x_i) + L(p_0|y_1)P(y_1|x_i),$$

gdzie

- $P(p_1|x_i)$, $P(p_0|x_i)$ - oznaczają estymowane przez model prawdopodobieństwo i oczywiście $P(p_0|x) = 1 - P(p_1|x)$,
- $L(p_i|y_j)$ oraz $i, j \in \{l, f\}$ - funkcja kosztu,
- x_i - i-ta obserwacja.

Klasyfikacja danej obserwacji jest opisana następującą nierównością:

$$R(p_1|x_i) \leq R(p_0|x_i)$$

Ona oznacza, że klasyfikujemy dany przykład jako pozytywny, jeżeli ryzyko takiej decyzji jest mniejsze niż zaklasyfikowania danej obserwacji jako negatywnej. Po przeprowadzaniu odpowiednich przekształceń algebraicznych otrzymujemy następujący wzór na klasyfikację przykładu jako pozytywny:

$$P(p_1|x_i) \geq \frac{L(p_1|y_0) - L(p_0|y_0)}{L(p_0|y_1) - L(p_1|y_1) - L(p_0|y_0) + L(p_1|y_0)}.$$

W przypadku gdy za funkcję kosztu przyjmiemy odpowiednią wartość z macierzy kosztu, to otrzymujemy następujący prób decyzyjny:

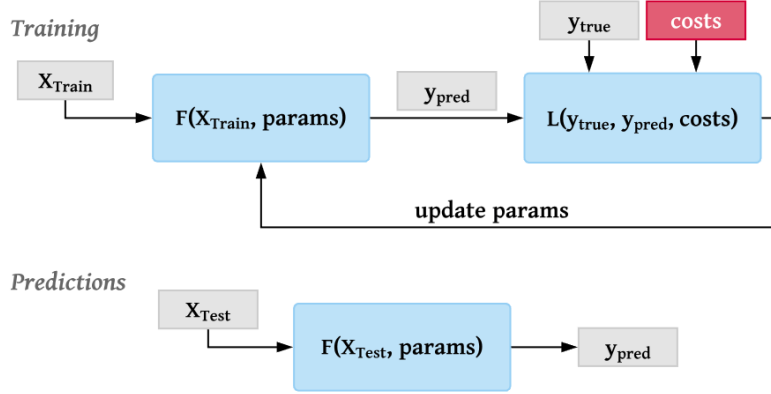
$$p \geq \frac{C_{FP_i} - C_{TN_i}}{C_{FN_i} - C_{TP_i} - C_{TN_i} + C_{FP_i}}.$$

2.4 Cost Sensitive Training

Drugą podgrupą metod wrażliwych na koszt jest *Cost Sensitive Trainig*. Są to metody, które już na etapie treningu modelu biorą pod uwagę koszt związany z klasyfikacją danej obserwacji. Schemat uczenia modelu jest przedstawiony na Rysunku 2.3. Model jako wejście otrzymuje zbiór danych treningowych, następnie dokonuje predykcji i na bazie prawdziwych odpowiedzi oraz kosztów wyznaczana jest skuteczność modelu i kolejno aktualizowane są wagi, a cały proces jest iteracyjnie powtarzany aż do momentu osiągnięcia zadanego kryterium stopu.

2.4.1 Regresja logistyczna wrażliwa na koszt

Pierwszym modelem, który stosunkowo łatwo przystosować do bycia wrażliwym na koszt jest regresja logistyczna. Kontynuując rozważania z podrozdziału 2.2.1



Rysunek 2.3: Schemat przedstawiający proces uczenia modelu wrażliwego na koszt. Źródło: <https://towardsdatascience.com/fraud-detection-with-cost-sensitive-machine-learning-24b8760d35d9>

chcemy, aby funkcja straty przyjmowała następujące wartości, które odpowiadają wartościom z macierzy kosztu:

$$J_i^c(\theta) = \begin{cases} C_{TP_i}, & \text{if } y_i = 1 \text{ and } h_\theta(\mathbf{x}_i) \approx 1, \\ C_{TN_i}, & \text{if } y_i = 0 \text{ and } h_\theta(\mathbf{x}_i) \approx 0, \\ C_{FP_i}, & \text{if } y_i = 0 \text{ and } h_\theta(\mathbf{x}_i) \approx 1, \\ C_{FN_i}, & \text{if } y_i = 1 \text{ and } h_\theta(\mathbf{x}_i) \approx 0. \end{cases}$$

W rezultacie funkcją, która zachowuje się w powyższy sposób jest:

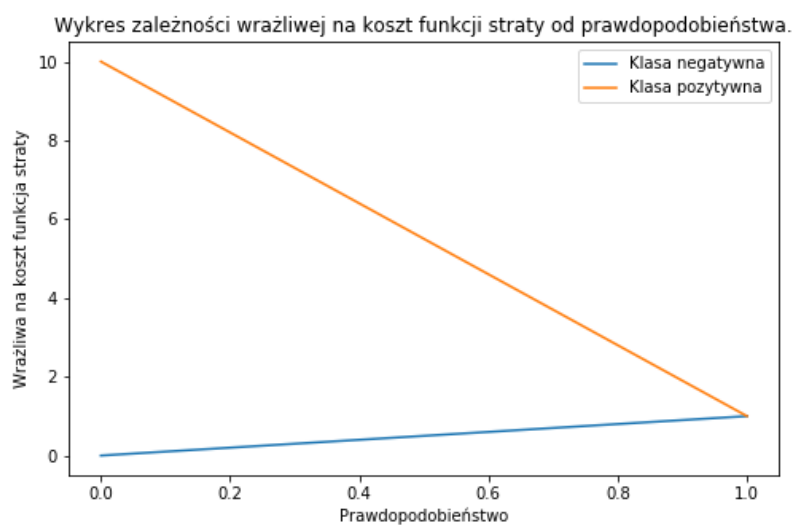
$$J^c(\theta) = \frac{1}{N} \sum_{i=1}^N \left(y_i \left(h_\theta(\mathbf{x}_i) C_{TP_i} + (1 - h_\theta(\mathbf{x}_i)) C_{FN_i} \right) + (1 - y_i) \left(h_\theta(\mathbf{x}_i) C_{FP_i} + (1 - h_\theta(\mathbf{x}_i)) C_{TN_i} \right) \right)$$

Następnie standardowo wykorzystując algorytm optymalizujący, który znajdzie odpowiednie współczynniki modelu trenujemy model predykcyjny.

2.4.2 Drzewo decyzyjne wrażliwe na koszt

Analogicznie jak w przypadku regresji logistycznej w celu wprowadzenia kosztu do treningu drzewa decyzyjnego musimy uzależnić proces powstawania modelu od kosztu danej próbki. Naturalnym miejscem dla drzewa decyzyjnego jest proces podziału danego węzła na kolejne węzły bądź liście. Dlatego definiujemy następującą miarę zanieczyszczenia, która następujący wzór:

$$I_c(\mathcal{S}) = \min \{ \text{Cost}(f_0(\mathcal{S})), \text{Cost}(f_1(\mathcal{S})) \}.$$



Rysunek 2.4: Wykres zależności wrażliwej na koszt funkcji straty od prawdopodobieństwa predykcji danej obserwacji. Kolor niebieski przedstawia wykres dla próbki o prawdziwej klasie pozytywnej, natomiast kolor pomarańczowy dla negatywnej. Wykres dla przykładowych wartości: $C_{TP} = 1$, $C_{FN} = 10$, $C_{FP} = 1$, $C_{TN} = 0$.

Dodatkowo w przy dokonywaniu predykcji klasyfikacja na wartość pozytywną bądź negatywną następuje wg następującego kryterium:

$$f(\mathcal{S}) = \begin{cases} 0, & \text{jeżeli } \text{Cost}(f_0(\mathcal{S})) \leq \text{Cost}(f_1(\mathcal{S})), \\ 1, & \text{w przeciwnym wypadku,} \end{cases}$$

gdzie jak zawsze 1 oznacza wartość pozytywną, a 0 negatywną. Z poprzednich rozważań oczywiście możliwe jest stworzenie modeli typu *ensemble*, w których klasyfikatorem bazowym byłoby drzewo decyzyjne wrażliwe na koszt, natomiast jest to poza obszarem badań aktualnej pracy.

Rozdział 3

Eksperyment

W celu porównania skuteczności metod wrażliwych na koszt względem siebie wykonamy eksperyment, którego celem będzie sprawdzenie, który z modeli lub ich kombinacja przyniesie największe oszczędności na zbiorze danych dot. oszustw z wykorzystaniem kart kredytowych. Jako dodatkowy wskaźnik wykorzystamy miarę skuteczności modelu F1 Score, która będzie w stanie odpowiedzieć jak zmieniła się ogólna skuteczność modelu.

3.1 Dane

Do eksperymentu zostanie wykorzystany zbiór danych Credit Card Fraud Detection zawierający 284,807 transakcji w tym zaledwie 492 oszustw. Tabela składa się z 30 kolumn, w tym 28 z nich są to zanonimizowane zmienne numeryczne, które były wcześniej poddane transformacji PCA (*ang. Principal Component Analysis*), a dodatkowo posiadamy informacje dot. czasu transakcji oraz kwoty. Ponadto wśród danych nie ma brakujących wartości. Podczas modelowania pominiemy zmienną czasową, ponieważ sama w sobie nie zawiera ona istotnej informacji, natomiast stworzenie znaczących zmiennych nie jest istotą przeprowadzanego eksperymentu.

Mimo tego, że dane zostały poddane anonimizacji, to na podstawie ?? możemy domyślać się z jakimi zmiennymi mieliśmy do czynienia przed transformacją. Podczas procesu dokonywania transakcji standardowo zbierane są: ID klienta, data dokonania transakcji, kwota, lokalizacja, typ transakcji (przykładowo płatność internetowa, płatność stacjonarna, wypłata z bankomatu), beneficjent transakcji (przykładowo linie lotnicze, hotel, wypożyczalnia samochodów itp.), historyczna informacja dot. czy dana transakcja była oszustwem, wiek klienta, kraj zamieszkania, kod pocztowy, typ karty (przykładowo VISA, MasterCard itp.). Na podstawie wymienionych informacji tworzy się agregaty czasowe bazujące na historii, które składają się z następujących kombinacji zmiennych:

- Klient, karta kredytowa;

- Typ transakcji, beneficjent transakcji, kategoria beneficjenta, kraj;
- W ciągu ostatnich godzin/dni/tygodni/miesięcy;
- Ilość, średnia lub suma transakcji.

Przykładowymi zmiennymi, które powstają w tym procesie są: Ilość transakcji dla tego samego klienta w ciągu ostatnich 6/12/24 godzin, Średnia wartość transakcji dla danego klienta z ostatniego tygodnia.

Informacja, która jest najbardziej istotna z punktu widzenia naszego modelowania, to rozkład oraz charakterystyka zmiennej dot. kwoty transakcji. Na Rysunku ?? znajduje się wykres rozkładu ...

Warto zauważyć, że najpopularniejszymi wartościami transakcji są:

Na podstawie przeprowadzonej analizy definiujemy w Tabeli 3.1 macierz kosztu dla tego eksperymentu, gdzie C_a to koszt administracyjny podjęcia sprawy przez analityka, który sprawdza czy dana obserwacja jest faktycznie oszustwem niezależnie od tego jaki był jego ostateczny werdykt, Amt_i to wartość transakcji jaką utracimy w przypadku nie wskazanie danej obserwacji jako oszustwa, natomiast zerowa wartość kosztu dla normalnej transakcji, która była prawidłowo wskazana, wynika z braku konieczności podjęcia inwestycji oraz strat z tego wynikających.

		Predykcja	
		Oszustwo	Normalna
Prawda	Oszustwo	C_a	Amt_i
	Normalna	C_a	0

Tabela 3.1: Macierz kosztu dla eksperymentu.

3.2 Opis eksperymentu

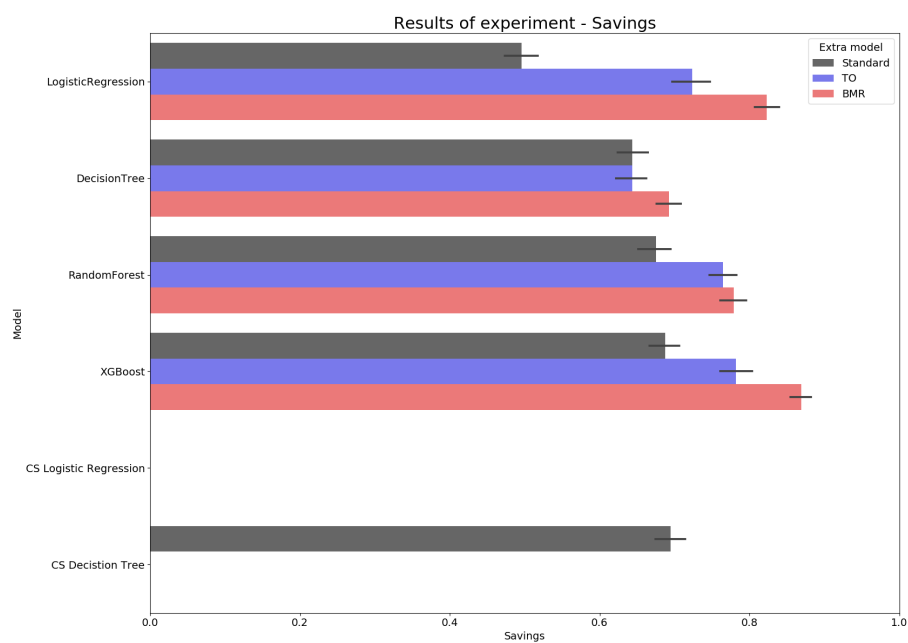
Eksperyment został przeprowadzony w następujący sposób: 50-krotnie dzielimy zbiór danych w proporcjach 50:17:33 na zbiór treningowy, walidacyjny oraz testowy. Następnie uczymy wszystkie modele na zbiorze treningowym. Dla modelu XGBoost wykorzystujemy zbiór walidacyjny do procesu wczesnego zatrzymywania (*ang. Early stopping*), natomiast dla modeli BMR oraz TO korzystamy z tego zbioru jako zbiór treningowy. Następnie dla wszystkich modeli dokonujemy predykcji na zbiorze testowym i mierzymy skuteczność typowań. Warto wspomnieć, że 50-krotne losowanie nowych podziałów zbioru ma na celu zmniejszenie ryzyka, że wynik zależy od wylosowanej próbki. Ponadto wykorzystano tę technikę w opozycji do standardowej warstwowej walidacji krzyżowej (*ang. Stratified Cross Validation*) z uwagi na chęć uniknięcia zbyt małej próbki testowej w przypadku walidacji krzyżowej z dużą ilością podziałów. Następnie dla

każdego modelu została obliczona wartość średnia oraz odchylenie standardowe dla wcześniej określonych miar.

Do implementacji skryptów został wykorzystany język programowania Python wraz z bibliotekami costcla, sklearn, pandas, numpy, matplotlib oraz język programowania bash.

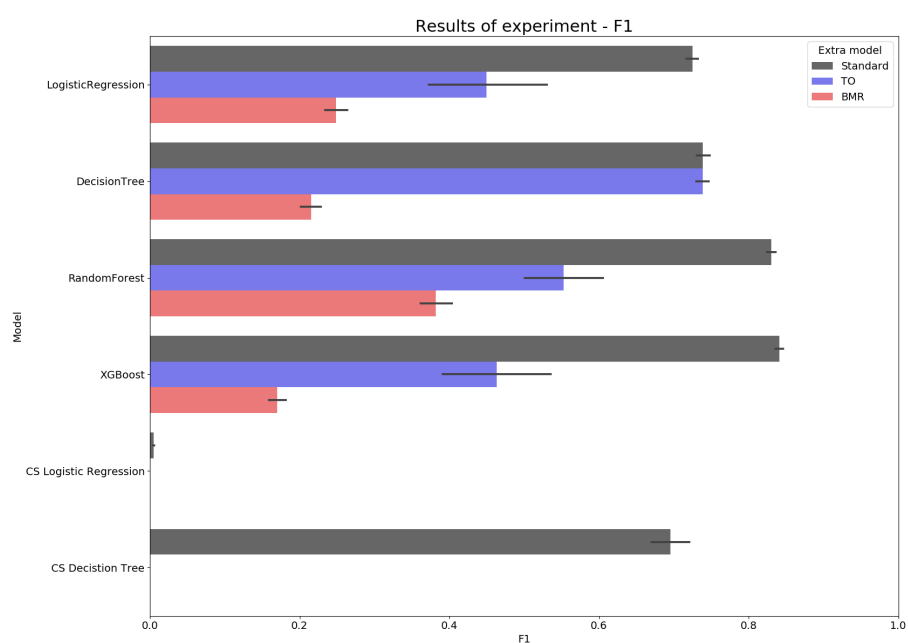
3.3 Wyniki

Lorem ipsum ...



Rysunek 3.1: Wyniki eksperymentu dla miary skuteczności Oszczędności.

Korzystając z Rysunku 3.3 możemy zauważyć, że: ...
Korzystając z Rysunku 3.3 możemy zauważyć, że: ...



Rysunek 3.2: Wyniki eksperymentu dla miary skuteczności F1-Score.

Rozdział 4

Podsumowanie