



Studium Licencjackie

Kierunek: Metody Ilościowe w Ekonomii i Systemy Informacyjne

Imię i nazwisko autora:

Jan Jarco

Nr albumu: 82762

Wykorzystanie metod uczenia maszynowego w wykrywaniu oszustw reklamowych

Praca licencjacka

pod kierunkiem naukowym

dr Małgorzaty Wrzosek

Instytut Ekonometrii

Warszawa 2021

Spis treści

Wprowadzenie	5
1 Reklama cyfrowa	7
1.1 Znaczenie reklamy cyfrowej	7
1.2 Rodzaje reklamy w marketingu cyfrowym	8
1.3 Reklama cyfrowa a handel elektroniczny	9
2 Mierzenie efektów reklamy cyfrowej	11
2.1 Sposób śledzenia aktywności użytkowników w Internecie	11
2.2 Model dochodowy marketingu cyfrowego	12
2.3 Miary efektów kampanii	13
3 Oszustwa w marketingu cyfrowym	15
3.1 Wpływ oszustw reklamowych na wyniki kampanii	15
3.2 Wpływ oszustw na przykładzie fałszywych kliknięć	16
4 Zastosowanie algorytmów uczenia maszynowego do predykcji fałszywych kliknięć	19
4.1 Opis eksperymentu i danych użytych do analizy	19
4.2 Opis działania algorytmów użytych do predykcji	20
4.3 Wyznaczanie optymalnych hiperparametrów	22
5 Analiza efektywności klasyfikacji na przykładzie modelu LightGBM	24
5.1 Zbiór danych użyty do klasyfikacji	24
5.2 Wstępna analiza danych	26
5.3 Optymalizacja hiperparametrów modeli klasyfikacyjnych	28
5.4 Struktura modelu	29
5.5 Wyniki klasyfikacji	30
5.6 Walidacja modelu	31
5.7 Przykładowe zastosowanie modelu w praktyce biznesowej	34
6 Wnioski	36
Bibliografia	37

Spis rysunków	39
Spis tabel	40
Streszczenie.....	41

Wprowadzenie

Rozwój usług handlowych, medycznych czy rozrywkowych ulokowanych w internecie postępuje w bardzo szybkim tempie, został on przyspieszony dodatkowo przez pandemię COVID-19, która w wyniku zamykania stacjonarnych punktów usługowych zmusiła usługodawców do zmiany sposobu prowadzenia swoich usług. Rozwój ten nie mógłby postępować bez odpowiednich narzędzi do informowania potencjalnych klientów o dostępnej ofercie usługowej, w tym przypadku najskuteczniejszym jest reklama cyfrowa. Jak pokazują badania, ten rodzaj reklamy każdego roku ma coraz większy udział w rynku, w zeszłym roku po raz pierwszy przekroczył 50% udziału w rynku reklamowym na świecie¹.

Ze względu na szybki rozwój usług marketingowych w internecie, reklamodawcy są zmuszeni mierzyć się z coraz to nowymi zagrożeniami w postaci oszustw reklamowych, które przybierają różne formy. Skala takich działań jest ogromna, szacowana na 42 mld dol. amerykańskich, według raportu eMarketer stawia to je wśród najbardziej dochodowych nielegalnych działalności². W wyniku tego procederu reklamodawcy tracą znaczące części budżetu na rzecz podmiotów niedziałających zgodnie z prawem. Dodatkowo statystyki powstałe na skutek mierzenia efektów kampanii marketingowych stają się mało wiarygodne i uniemożliwiają dostosowanie targetowania reklamy, co prowadzi do utraty potencjalnych zysków. Szczególnym rodzajem oszustw reklamowych jest zjawisko fałszywych kliknięć, czyli generowania nadmiarowej liczby kliknięć w reklamę rozliczaną w systemie *Pay-per-click*, które nie mają przełożenia na poprawienie wyników sprzedażowych. W celu ograniczenia strat wynikłych z tego zjawiska branża marketingowa uruchamia usługi, które mają za zadanie wykrywanie fałszywego ruchu na stronie internetowej. Narzędzia te oparte są na modelach uczenia maszynowego i mają bardzo wysoką skuteczność w przewidywaniu poprawnego ruchu internetowego, który odznacza się pożądaną akcją użytkownika, określaną dalej jako konwersję w postaci dokonania zakupu, dodania produktu do koszyka zakupowego, pobrania aplikacji czy przeglądania strony internetowej przez określony czas.

Celem niniejszej pracy jest uzyskanie narzędzia wspomagającego ocenę, czy dany użytkownik jest oszustem reklamowym. Do budowy takiego narzędzia zostaną wykorzystane dane od chińskiego lidera usług marketingowych na urządzenia mobilne dotyczące dokonania

¹ Cramer-Flood (2020).

² Perrin (2020).

konwersji, w tym przypadku pobrania aplikacji, po pojedynczych kliknięciach użytkowników w reklamę. Na podstawie tych danych, dotyczących m.in. powtarzalności kliknięć w czasie, pory dnia, kanału reklamowego i wersji oprogramowania urządzenia skonstruowany zostanie model klasyfikujący, czy zostanie wykonana konwersja po danych kliknięciach. W celu identyfikacji potencjalnych oszustów reklamowych zostanie opracowana specjalna reguła decyzyjna, wykorzystująca wyniki z wspomnianego modelu klasyfikacyjnego w postaci średnich wartości prawdopodobieństwa wykonania konwersji po kliknięciach w reklamę, dla danego użytkownika oraz ilość kliknięć użytkownika w reklamę w czasie.

W rozdziale pierwszym pracy została opisana sytuacja na rynku reklamy internetowej wraz z jej rosnącym znaczeniem względem reklamy tradycyjnej, zostały także przedstawione rodzaje reklamy internetowej wraz z jej wpływem na handel elektroniczny. Rozdział drugi opisuje sposób śledzenia aktywności użytkowników na stronach internetowych i jego wykorzystanie w ramach rozliczania kosztów za reklamę w internecie. Zjawisko oszustw reklamowych i jego negatywne skutki, z którymi muszą się zmagać reklamodawcy zostały scharakteryzowane w rozdziale trzecim, tu także wyszczególnione zostało zjawisko fałszywych kliknięć i możliwe sposoby ograniczenia strat wynikłych z tej nieuczciwej praktyki. Rozdział czwarty prezentuje zastosowane podejście oraz proponowane metody wspomagające ograniczenie strat wynikłych z problemu zjawiska fałszywych kliknięć. Dokładny opis użytego zbioru danych, analizę efektywności klasyfikacji modelu oraz przykładową implementację reguły klasyfikującej potencjalnych oszustów reklamowych zawarto w rozdziale piątym. Rozdział szósty stanowi podsumowanie niniejszej pracy, przedstawia główne wnioski i możliwości poprawy narzędzia ograniczającego negatywne skutki zjawiska fałszywych kliknięć.

1 Reklama cyfrowa

1.1 Znaczenie reklamy cyfrowej

Rzeczywistość cyfrowa różni się w sposób znaczący od tej tradycyjnej. Wraz ze zwiększoną aktywnością konsumentów w internecie i rozwojem usług internetowych rośnie rynek reklamy cyfrowej, która ma największą skuteczność w docieraniu do potencjalnego klienta m.in. w handlu elektronicznym. Jako główne przyczyny tej przewagi wyróżnia się następujące czynniki:

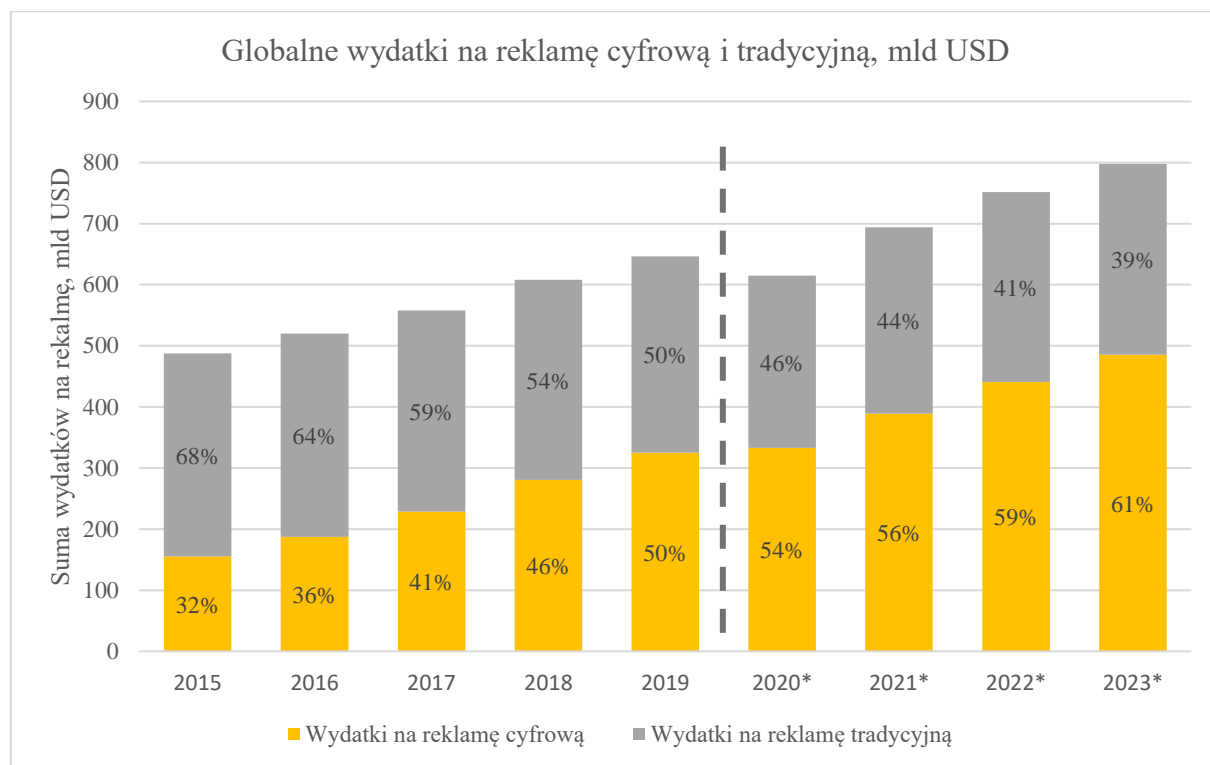
1. Usługi marketingu cyfrowego pozwalają na dotarcie do potencjalnego klienta lokalnie jak i na całym świecie, dzięki bardzo małym ograniczeniom w dostępie do informacji.
2. Natychmiastowy i często nieograniczony czasowo dostęp do informacji.
3. Rezultaty reklamy mogą być stosunkowo łatwe do zmierzenia względem reklamy tradycyjnej.
4. Możliwość budowania bezpośredniej relacji dzięki reklamie w mediach społecznościowych.

Z powodu wyżej wymienionych zalet reklamy cyfrowej względem tradycyjnej każdego roku rośnie jej udział w wydatkach na reklamę na całym świecie. Rysunek 1 przedstawia globalne wydatki na reklamę w podziale na reklamę cyfrową i tradycyjną. Pokazuje on, że poziom wydatków na reklamę cyfrową w 2019 roku osiągnął po raz pierwszy poziom powyżej 50% globalnych wydatków na reklamę. Według prognoz agencji eMarketer w ciągu najbliższych trzech lat udział marketingu cyfrowego wzrośnie do 61%, z każdym rokiem wypierając marketing tradycyjny³.

Należy jednak zauważyć, że reklama tradycyjna ma wciąż duże znaczenie, gdyż jest ona dużo skuteczniejszą formą dotarcia do osób starszych, które w mniejszym stopniu korzystają z usług internetowych, dlatego marketing cyfrowy jest często nakierowany na młodszą grupę odbiorców, a tradycyjny celuje w osoby z starszej grupy wiekowej. Odpowiednie prowadzenie tych dwóch form marketingu jednocześnie powoduje, na bazie efektu synergii, dodatkowe zwiększenia zysków dla przedsiębiorstwa, m.in. dzięki dostosowywaniu reklamy do różnych grup wiekowych.⁴

³ Cramer-Flood (2020, s. 4).

⁴ Jothi (2019, s. 5).



Rysunek 1: Globalny udział wydatków na reklamę cyfrową mld \$. Wykres przedstawia globalne wydatki na reklamę online i offline według raportu eMarketer. Wykres pokazuje dane z lat 2015 – 2019 i prognozy na lata 2020 – 2023. Na podstawie wykresu można zauważyć, że obserwujemy stały wzrost udziału wydatków na reklamę online, który osiągnął 50% w roku 2019 (Źródło: Cramer-Flood, 2020).

1.2 Rodzaje reklamy w marketingu cyfrowym

Rozwój rynku reklamy w internecie ma swoje odzwierciedlenie także w różnorodności tych usług. Wraz z ewolucją środowiska internetowego powstają naturalnie nowe możliwości i rodzaje reklamy⁵. Do podstawowych narzędzi należą wymienione poniżej:

1. Marketing e-mailowy (ang. *E-mail marketing*)

Reklama jest indywidualnie dystrybuowana do użytkowników przy pomocy wiadomości email. Ten sposób daje duże możliwości personalizowania zawartości na podstawie preferencji odbiorcy. Dzięki efektywnemu wykorzystaniu oprogramowania istnieje możliwość trafnego segmentowania grup odbiorców i trafiania z przekazem w odpowiednim czasie dostosowanym do celów marketingowych

2. Marketing w sieciach społecznościowych (ang. *social media marketing*)

⁵ Jothi (2019, s. 3–4).

Liczne serwisy społecznościowe, takie jak Facebook, Instagram, LinkedIn czy YouTube, wykorzystują swoje grono aktywnych użytkowników do wyświetlania spersonalizowanych reklam w ich kanałach aktualności. Takie rozwiązanie jest nakierowane na korzyści wynikające z możliwości udostępniania, komentowania czy dodawania reakcji pod postami reklamodawców, co jest widoczne dla innych użytkowników sieci, głównie osób powiązanych z kontem użytkownika.

3. Marketing w wyszukiwarkach internetowych (ang. *search engine marketing*)

Dzięki użyciu narzędzi w postaci wyszukiwarek internetowych jak Google czy Bing reklama wyświetlana jest w momencie wyszukiwania przez użytkownika odpowiednich haseł w wyszukiwarce. Odpłatne linki przekierowują bezpośrednio na stronę internetową reklamodawcy zwiększając przy tym ruch na stronie internetowej. Zaawansowane techniki optymalizacji wyszukiwań pozwalają na inteligentne dopasowanie treści do przedmiotu zainteresowania szukającego.

4. Reklama na urządzeniach mobilnych (ang. *mobile marketing*)

Za sprawą rosnącej popularności urządzeń mobilnych (smartfonów, tabletów), na które można pobierać liczne aplikacje oferujące różne narzędzia, reklamodawcy umieszczają wewnątrz aplikacji reklamy, rozliczając się z właścicielami aplikacji za udostępnianie przestrzeni reklamowej. Do użytkownika można dotrzeć za sprawą krótkich powiadomień generowanych przez aplikacje, reklam w wyszukiwarkach urządzeń mobilnych, czy reklam wyświetlanych wewnątrz aplikacji w czasie korzystania.

5. Reklama wyświetlana na stronach internetowych (ang. *Digital display marketing*)

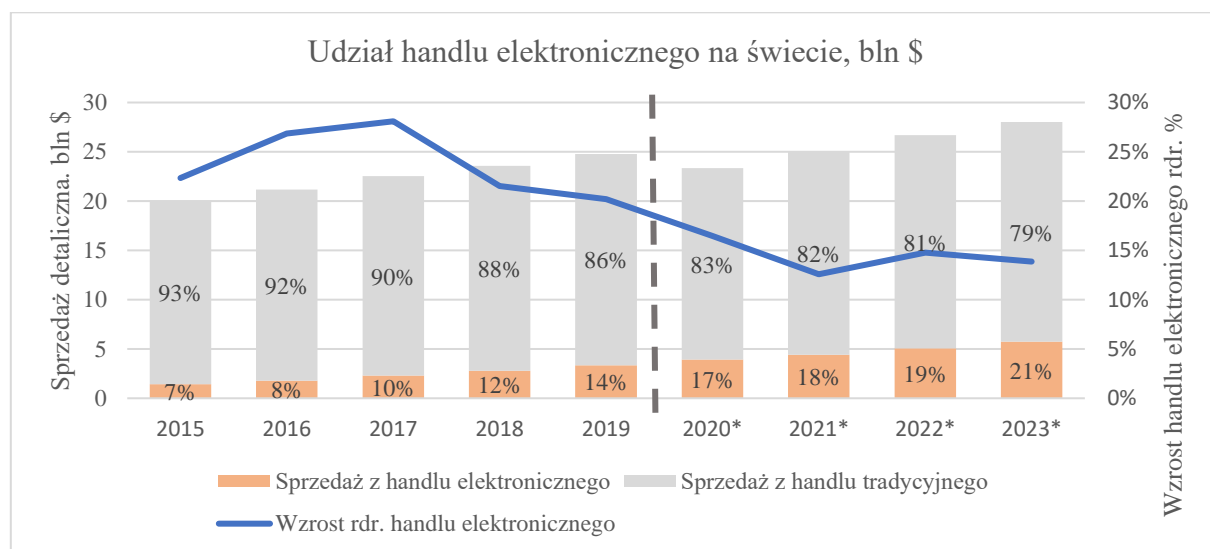
Właściciele stron internetowych udostępniają reklamodawcom przestrzeń na swoich stronach i otrzymują zapłatę za liczbę kliknięć w wyświetloną reklamę.

Oszustwa w reklamie cyfrowej występują w każdym z jej rodzajów, jednak ich skala jest różna. Najbardziej zagrożone są: reklama w wyszukiwarkach internetowych, na urządzeniach mobilnych i wyświetlana na stronach internetowych.

1.3 Reklama cyfrowa a handel elektroniczny

Dzięki coraz większemu dostępowi do usług internetowych rośnie znaczenie handlu elektronicznego wśród całości sprzedaży detalicznej. Według szacunków, w latach 2020 – 2024, wzrost e-handlu będzie postępował ze skumulowanym współczynnikiem wskaźniku

wzrostu na poziomie 13.7%. Taka sytuacja jest okazją dla sprzedawców do zwiększenia liczby swoich potencjalnych klientów poprzez wykorzystanie internetu jako platformy do zawierania transakcji. Handel elektroniczny posiada wiele zalet, które pozwalają na zwiększenie zysków przez handlowców.



Rysunek 2: Globalny udział handlu elektronicznego i jego roczny wzrost w latach 2015-2023, dane od roku 2020 są wartościami prognozowanymi (Źródło: Cramer-Flood, 2020b).

Jak pokazuje Rysunek 2, sprzedaż elektroniczna rośnie od kilku lat. Według raportu pandemia wywołana przez koronawirusa w pierwszym kwartale 2020 roku zauważalnie zwiększyła udział handlu elektronicznego. Między rokiem 2019 i 2020 jego udział wzrósł z 14% do 17% względem całości sprzedaży detalicznej⁶. Ten szybki wzrost był wywołany ograniczeniem w liczbie otwartych sklepów tradycyjnych, przez co klienci często byli zmuszeni zdecydować się na zakupy w internecie, które niosą mniejsze ryzyko zakażenia.

Najsukuteczniejszym sposobem reklamowania się przez e-handlowców są właśnie cyfrowe kanały informacyjne. Mają największą skuteczność dzięki łatwości przekierowania potencjalnego konsumenta natychmiast na platformę sprzedażową, gdzie dane usługi lub towary mogą zostać nabyte przez konsumenta, z przesyłką bezpośrednio do miejsca zamieszkania. Aby ocenić skuteczność reklamy cyfrowej niezbędna jest możliwość jej precyzyjnego mierzenia i oceniania, w celu optymalizacji działań reklamowych. W tym celu skonstruowano liczne miary skuteczności, które zostały opisane w następnym rozdziale niniejszej pracy.

⁶ Sirimanne, S. (2021). COVID-19 and e-commerce: a global review. Pobrane z: <https://unctad.org/webflyer/covid-19-and-e-commerce-global-review>. (Data dostępu: 22 marca 2021r.).

2 Mierzenie efektów reklamy cyfrowej

2.1 Sposób śledzenia aktywności użytkowników w Internecie

W trakcie przeglądania internetu użytkownicy poprzez swoją aktywność zostawiają bardzo dużo informacji, podstawowych takich jak: typ urządzenia, wersja przeglądarki, indywidualny numer IP, czas i lokalizacja aktywności oraz bardziej szczegółowych, zależnych od sposobu rejestracji aktywności, takich jak zachowania, czy preferencje użytkownika.

Głównym powodem, dla którego dane o przeglądaniu stron internetowych są zbierane przez właścicieli stron i reklamodawców jest zdobycie potrzebnej wiedzy na temat użytkowników – ich zachowań oraz preferencji. Ta wiedza, budowana dzięki danym o ruchu na stronie internetowej służy później do spersonalizowania rekomendacji produktów lub reklam, które mogłyby spodobać się użytkownikowi i zachęciłyby do zakupu usługi lub produktu. Drugim powodem, który stoi za koniecznością rejestrowania aktywności użytkowników na stronie internetowej jest optymalizacja wyglądu strony internetowej, dla lepszego doświadczenia użytkownika z jej korzystania. Przejrzysty układ i wygląd strony internetowej powodują większe zadowolenie klientów z korzystania z niej. Pozytywne doświadczenia użytkownika powodują większe przychody z sprzedaży na danej platformie internetowej⁷. Kolejnym celem, który jest osiągany dzięki śledzeniu aktywności jest tzw. „retargeting”. Po odwiedzeniu witryny sklepu internetowego z danym produktem użytkownikowi pokazują się reklamy przedstawiające dany artykuł, które przekierowują bezpośrednio na stronę reklamodawcy. Mają one na celu przypomnieć użytkownikowi o ofercie sklepu, którą wcześniej przeglądał.

Reklamodawcy korzystają ze śledzenia działań klienta z użyciem narzędzi właściciela strony lub podmiotów trzecich. Najczęściej opierają się one na wykorzystaniu tzw. „plików cookie”. Są to dane informatyczne o małym rozmiarze, przechowywane na urządzeniu użytkownika, generowane każdorazowo za jego zgodą, przy pomocy zatwierdzenia komunikatu pokazywanego podczas pierwszego otwarcia strony.⁸ Właściciele strony lub podmioty trzecie, które generują te pliki mają stały dostęp do tych plików i robią następne zapisy informacji wewnątrz nich. Zależnie od typu pliki cookie mają swój okres żywotności,

⁷ Crawford, E. (2020). Website Tracking: Why and How Do Websites Track You? Pobrane z: <https://www.cookiepro.com/blog/website-tracking/>. (Data dostępu: 27 marca 2021r.).

⁸ home.pl. (2020). Czym są pliki cookies (ciasteczka) w przeglądarce? Pobrane z: <https://pomoc.home.pl/baza-wiedzy/czym-wlasciwie-sa-pliki-cookies-ciasteczka-w-przegladarce>. (Data dostępu: 27 marca 2021r.).

który zazwyczaj wynosi 30 dni od powstania. Okres ten może mieć inną długość w zależności od ustawień właściciela witryny.

Pliki cookie są dzielone według kryterium czasu żywotności na:

- pliki sesyjne – śledzą użytkownika wyłącznie podczas przeglądania strony
- pliki permanentne – pozostają na bardzo długi czas na urządzeniu użytkownika, dopóki nie zostaną usunięte z pamięci przeglądarki automatycznie przez przeglądarkę lub przez decyzję użytkownika

Takie dane pozyskane z narzędzi do śledzenia posłużą w dalszych rozdziałach pracy do klasyfikacji kliknięć, po których istnieje większe prawdopodobieństwo wykonania konwersji, w postaci pobrania aplikacji na urządzenie.

2.2 Model dochodowy marketingu cyfrowego

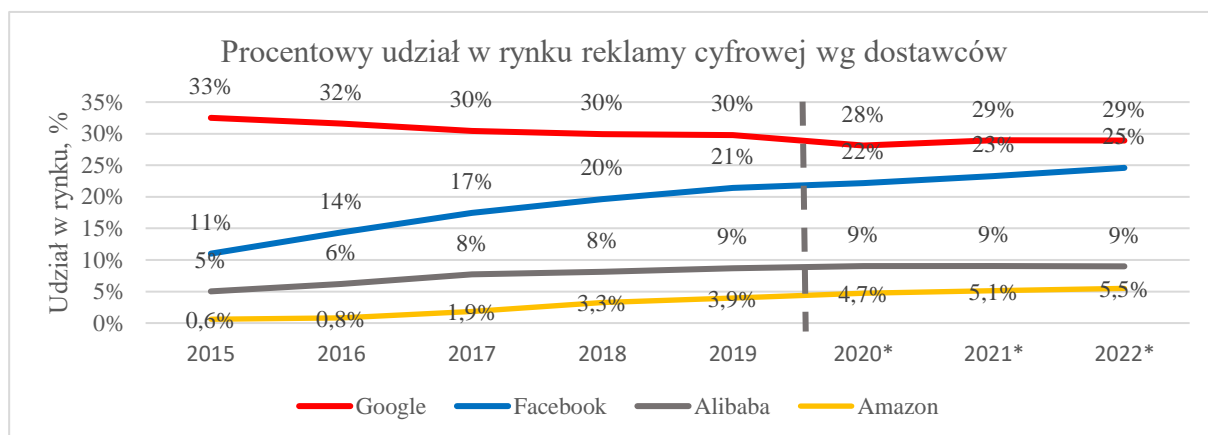
Istnieją dwa główne modele rozliczania się dostawców reklamy internetowej z reklamodawcami:

- płatność za kliknięcia w reklamę (ang. *Pay-per-click*),
- płatność za wyświetlenia reklamy, inaczej impresje (ang. *Pay-per-view*).

Aktualnie bardziej znaczącym modelem rozliczania się w marketingu cyfrowym jest system płatności za kliknięcie (*Pay-per-click*). W tym modelu reklamodawcy płacą ustaloną kwotę dostawcy usług reklamowych za każde kliknięcie w ich reklamę wyświetlaną użytkownikowi. Jest on wykorzystywany najczęściej przy reklamie wyświetlanej na stronach internetowych lub w wyszukiwarkach internetowych w trakcie pozycjonowania odpowiednich linków.

Co obrazuje Rysunek 3, przedstawiający udział poszczególnych dostawców w rynku reklamy internetowej, największymi jej dostawcami są firmy Google, Facebook, Alibaba i Amazon. Te cztery firmy o dużym znaczeniu na rynku posługują się właśnie systemem *Pay-per-click* do rozliczania się z reklamodawcami. Widoczna jest duża stabilność udziałów dostawców posługujących się tym systemem w rynku reklamowym, co sprawia, że system ten ma aktualnie bardzo wysokie znaczenie i bardzo prawdopodobne, że będzie się ono utrzymywało w najbliższych latach.⁹

⁹ Cramer-Flood (2020, s. 8).



Rysunek 3: Najwięksi dostawcy usług reklamy cyfrowej, odpowiedzialni w roku 2020 za 64% światowych przychodów posługując się właśnie systemem Pay-per-click do rozliczania się ze swoimi klientami (Źródło: Cramer-Flood, 2020)

2.3 Miary efektów kampanii

Z opisanym systemem naliczania kosztów wiąże się bardzo istotna miara, czyli „koszt za kliknięcie”. Jest to metryka, przy pomocy której reklamodawcy rozliczają się z największymi dostawcami usług reklamy typu Pay-per-click. Wielkość tej metryki jest wyznaczana ze wzoru:

$$\text{Koszt za kliknięcie} = \frac{\text{Koszt reklamy}}{\text{Liczba kliknięć w reklamę}}$$

Koszt ten może być ustalany w formie stałej, niezmiennej lub w trakcie aukcji przeprowadzanych na odpowiednich platformach, gdzie reklamodawcy licytują pozycję swoich reklam ustalając maksymalny koszt za kliknięcie, który mogą zapłacić¹⁰. Przeciętny koszt jest zależny od dostawcy, branży, kraju, a także od ogólnego trendu rynkowego. W drugim kwartale 2020 r. przeciętny globalny koszt za kliknięcie wyniósł 0.49 dol. amerykańskiego¹¹.

Reklama ma na celu zachęcenie potencjalnych klientów do zakupu produktów i usług, dlatego istotną miarą dla marketerów jest współczynnik wejść na stronę internetową zakończonych późniejszym zakupem, bądź inną ważną dla działalności biznesowej aktywnością. Koszt reklamy musi być również analizowany równolegle ze współczynnikiem

¹⁰ WordStream. (2021). What Is PPC? Learn the Basics of Pay-Per-Click (PPC) Marketing. Pobrane z: <https://www.wordstream.com/ppc>. (Data dostępu: 02 kwietnia 2021r.).

¹¹ Statista Research Department. (2021). Search advertising cost-per-click (CPC) worldwide. Pobrane z: <https://www.statista.com/statistics/873639/search-advertising-cpc/>. (Data dostępu: 03 kwietnia 2021r.).

konwersji (z ang. *Conversion rate*), który wyraża udział osób, które dokonały konwersji na stronie po kontakcie z reklamą, w wartościach procentowych¹².

$$\text{Współczynnik konwersji (\%)} = \frac{\text{liczba konwersji}}{\text{liczba wizyt na stronie po kontakcie z reklamą}}$$

Niezależnie od branży, współczynnik ten nie przekracza 5%, chociaż dla reklam w wyszukiwarkach jest często kilkukrotnie wyższy niż w przypadku reklamy wyświetlanej na standardowej witrynie internetowej. Dzieje się tak, ponieważ użytkownik, który wyszukał odpowiednią frazę w wyszukiwarce internetowej, jest bardziej zainteresowany zakupem.

Powyższe miary są kluczowe dla działań marketingowych nakierowanych na sprzedaż produktów lub usług. Określają one w bezpośredni sposób sukces kampanii reklamowych w kształtowaniu sprzedaży. Opisane współczynniki określają poziom powodzenia działań w internecie, jednak mogą one być zniekształcane przez niechciany ruch w internecie.

¹² Google Ads. (2021). Conversion rate: Definition - Google Ads Help. Pobrane z: <https://support.google.com/google-ads/answer/2684489?hl=en>. (Data dostępu: 03 kwietnia 2021r.).

3 Oszustwa w marketingu cyfrowym

Wraz z szybkim rozwojem reklamy cyfrowej reklamodawcy ciągle zwiększają swoje budżety na właśnie ten typ reklamy. Nastąpiła też szybka automatyzacja procesów zakupu reklamy, jej publikacji i mierzenia efektów spowodowana głównie przez digitalizację tych procesów. Dzięki dogodnym, zautomatyzowanym narzędziom reklamodawcy mogą działać aktywnie na bardzo dużą skalę. Szybkie zwiększanie skali zakupów usług reklamowych w internecie spotęgowało zainteresowanie oszustów reklamowych, którzy za pomocą za pomocą swoich narzędzi sprawiają, że branża reklamowa bardzo dużo traci w wyniku nielegalnych wyłudzeń. Oszustwa reklamowe (ang. *ad fraud*) przybierają różne formy, właściwie na każdym etapie zakupu reklamy jest możliwość do prowadzenia nielegalnej działalności przynoszącej bardzo duże zyski. Oszustwa reklamowe definiuje się jako praktykę fałsyfikowania ruchu w sieci lub związanych z nim aktywności w celu naliczania opłat reklamodawcom za impresje (wyświetlenia reklamy), kliknięcia lub inne aktywności, które nie miały miejsca w rzeczywistości.

Jak dla każdej działalności nielegalnej zmierzenie skali zjawiska oszustw reklamowych jest bardzo trudne, dlatego szacunki sumy globalnych strat rynku reklamowego się bardzo różnią. Według agencji Juniper Research suma strat wyniosła 42 mld \$ w roku 2019¹³, jednocześnie w tym samym okresie GroupM, największa na świecie grupa agencji mediowych podaje, że branża traci 22.4 mld \$ rocznie¹⁴. Z pierwszego cytowanego raportu wynika, że ok. 13% budżetów przeznaczanych na reklamę cyfrową trafia do oszustów i nie przynosi zysków dla przedsiębiorstw. Niestety, szacuje się, że skala procederu będzie w kolejnych latach wzrastać i może osiągnąć nawet 100 mld \$ w roku 2023¹⁵.

3.1 Wpływ oszustw reklamowych na wyniki kampanii

Istnieje bardzo dużo różnych rodzajów działań nielegalnych, stosowanych przez oszustów, które polegają na manipulacji m.in. liczbą impresji, kliknięć, konwersji, instalacji aplikacji mobilnych, które wpływają bezpośrednio na miary opisane w rozdziale 2.3. W badaniu przeprowadzonym przez PPC Protect na 410 marketerach aż 31.8% respondentów wymienia oszustwa reklamowe (w tym fałszywe kliknięcia) wśród dwóch największych zagrożeń w marketingu cyfrowym¹⁶. Istnieją dwa główne powody, dla których reklamodawcy

¹³ Perrin (2020).

¹⁴ Cherico, M. (2019). Ad fraud costing \$22.4 billion globally says GroupM. Pobrane z: <https://www.foxbusiness.com/media/ad-fraud-22-billion-globally-groupm>. (Data dostępu: 23 marca 2021r.).

¹⁵ Perrin (2020).

¹⁶ PPC Protect. (2021). The Global Click Fraud Report 2020-2021. Pobrane z: <https://try.ppcprotect.com/click-fraud-report-2021/>. (Data dostępu: 03 kwietnia 2021r.).

bardzo poważnie traktują ten problem. Przede wszystkim jest to bezpośrednia strata budżetów reklamowych, które są przekazywane oszustom w wyniku kreowania fałszywego ruchu w sieci. Drugim powodem jest fałszowanie wyników kampanii, co sprawia, że działania analityczne nakierowane na optymalizację kierunkowania reklam cechują się mniejszą wiarygodnością.

3.2 Wpływ oszustw na przykładzie fałszywych kliknięć

Niniejsza praca opisuje przede wszystkim główny rodzaj działań fałszujących wyniki kampanii, czyli zjawisko fałszywych kliknięć (ang. *click fraud*). Ten proceder można zaliczyć do grupy działań mających na celu manipulowanie miarami ruchu w internecie, za które rozliczają się reklamodawcy. Do rozliczeń w systemie Pay-per-click wykorzystuje się liczbę m.in. kliknięć - marketerzy płacą za każde pojedyncze kliknięcie według stawki określanej jako „koszt za kliknięcie” (rozdział 2.3). System ten tworzy możliwość nadużyć w postaci nadmiarowych kliknięć w reklamę, które nie mają realnego przełożenia na zwiększenie wyników sprzedaży.

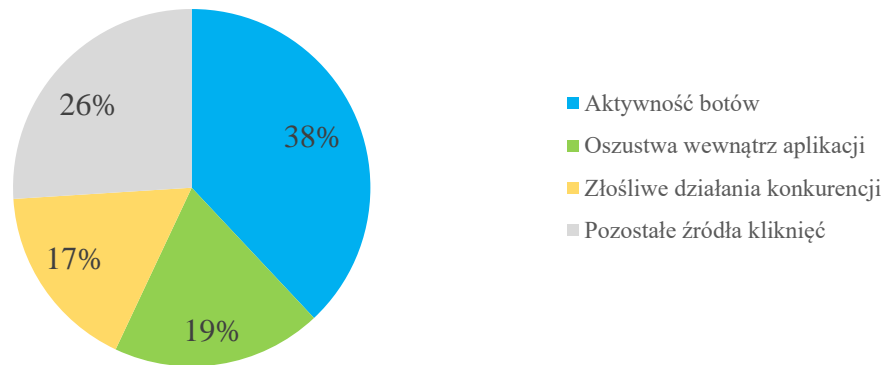
Największe sieci reklamowe takie jak Google Ads pełnią podwójną rolę – dają narzędzia do zwiększania zysków takie jak publikacja reklamy na stronach internetowych oraz umożliwiają bardziej korzystne pozycjonowanie linków w wynikach wyszukiwania. Taka działalność rodzi konflikt interesów, ponieważ tacy dostawcy z jednej strony dążą do optymalizacji kosztów reklamowych swoich klientów i jednocześnie mają na celu zwiększenie swoich zysków z publikacji reklam. Zwiększanie liczby kliknięć generuje dla dostawców korzyści, dlatego mogą potencjalnie za sprawą botów w sposób nieetyczny sztucznie zwiększać liczbę kliknięć i przy tym swoich zysków.¹⁷

Istnieją również podmioty trzecie, dla których nieefektywne wykorzystanie budżetów reklamowych przez przedsiębiorstwo generuje korzyści. Może to być nieuczciwa konkurencja reklamodawcy bądź konkurencja w postaci innych sieci reklamowych.

Dane wskazują, że aż 36% z wszystkich kliknięć w reklamy wyświetlane na stronach internetowych jest potencjalnie fałszywych. W przypadku reklamy w wyszukiwarkach internetowych ten udział wynosi dużo mniej – 11% wszystkich kliknięć, ze względu na mniejszą dochodowość tych działań dla zorganizowanych oszustów. W tym przypadku głównym źródłem nadużyć są jednak nieetyczne działania konkurencji.

¹⁷ Wilbur, Zhu (2008).

Procentowy udział źródeł fałszywych kliknięć według PPC Protect



Rysunek 4: Procentowy udział źródeł fałszywych kliknięć według PPC Protect pokazuje aktywność botów jako główne źródło problemu (Źródło: PPC Protect, 2021)

Wśród źródeł fałszywych kliknięć dominują rozwiązania wykorzystujące boty, czyli zautomatyzowane skrypty, które mają za zadanie masowo klikać w reklamy, generując przy tym większe koszty dla atakowanego reklamodawcy. Co ciekawe, według marketerów ponad połowa fałszywych kliknięć powstaje w wyniku działań nieuczciwej konkurencji, jednak w rzeczywistości udział ten jest znacznie niższy, wynosi 17%. Oprócz działań celowych nieprawidłowe kliknięcia mogą być również spowodowane przypadkowymi kliknięciami w reklamę przez zwykłych użytkowników (nie mających na celu generowania nadmiernych kosztów dla reklamodawcy), co tylko obciąża reklamodawcę, a nie generuje kolejnych konwersji. Jednak takie przypadkowe kliknięcia mogą też być wymuszane przez nieodpowiednie umieszczenie reklamy na stronie, na przykład takie jak nakładanie się wielu reklam na siebie. W tej sytuacji użytkownik klika wiele reklam jednocześnie zupełnie nie będąc tego świadomym.

Skala problemu fałszywego ruchu w internecie jest zróżnicowana w różnych regionach świata. Regiony EMEA i NAFTA wyróżniają się zdecydowanie niższym udziałem fałszywych kliknięć (poniżej 20%) względem regionu azjatyckiego, gdzie współczynnik ten wynosi średnio 24.2%¹⁸. Jednak niezależnie od tych różnic należy uznać, że skala problemu jest znacząca, a reklamodawcy ponoszą ogromne straty budżetowe. Z tego powodu powstają rozwiązania, które mają na celu zniwelować te straty przy pomocy zaawansowanych metod analitycznych. Celem branży marketingu cyfrowego jest osiągnięcie takiej sytuacji na rynku, że reklamodawcy będą

¹⁸ PPC Protect. (2021). The Global Click Fraud Report 2020-2021. Pobrane z: <https://try.ppcprotect.com/click-fraud-report-2021/>. (Data dostępu: 03 kwietnia 2021r.).

rozliczani wyłącznie z ruchu powstałego w wyniku prawdziwych działań, rzeczywistych użytkowników. W ciągu ostatniej dekady powstało wiele firm, których działania są nakierowane szczególnie na obniżenie strat wynikłych z fałszywych kliknięć (ClickCease, PPC Protect, Integrated Ad Science, White Ops i wiele innych). Ich rozwiązania polegają głównie na wykrywaniu podejrzanych pod względem jakości kliknięć przy pomocy wielu parametrów takich jak: powtarzalność kliknięć z tego samego adresu IP w czasie, podejrzane adresy IP, podejrzana lokalizacja, użycie usługi VPN, wyłączenie kodu JavaScript w przeglądarce i wiele innych.¹⁹ Działania te mają na celu m.in. wykrywanie fałszywego ruchu i blokowanie dostępu do stron urządzeniom należącym do nieuczciwych podmiotów. Dzięki temu reklamodawcy tracą mniejsze części budżetu²⁰.

W dalszej części pracy wykorzystane zostaną metody uczenia maszynowego w celu zbudowania algorytmu, który będzie przewidywał czy po kliknięciu w reklamę nastąpi pobranie aplikacji. Opis i porównanie modeli zostaną poprzedzone opisaniem algorytmów wykorzystanych do budowy reguły decyzyjnej.

¹⁹ Oentaryo i in. (2014).

²⁰ ClickCease. (2021). ClickCease™ - Click Fraud Protection & Prevention. Pobrane z: <https://www.clickcease.com/features.html>. (Data dostępu: 04 kwietnia 2021r.).

4 Zastosowanie algorytmów uczenia maszynowego do predykcji fałszywych kliknięć

Negatywne skutki oszustw w marketingu cyfrowym mogą być ograniczane przy pomocy odfiltrowywania ruchu internetowego pochodzącego od potencjalnych oszustów. Odfiltrowywanie na stronie ruchu, który nie przynosi korzyści dla reklamodawcy jest dużym wyzwaniem, ponieważ nie ma pewnej informacji, kto generuje kliknięcia w daną reklamę rozliczaną w systemie Pay-per-click. Istnieje wyłącznie możliwość oceny prawdopodobieństwa, czy dany użytkownik wykona w przyszłości pożądaną przez nas konwersję na podstawie jego wcześniejszej aktywności na stronie internetowej reklamodawcy. Oszacowane przy pomocy modeli klasyfikacyjnych uczenia maszynowego prawdopodobieństwa wykonania konwersji dla danego użytkownika mogą posłużyć jako wyznacznik oceny, czy jest on potencjalnym oszustem, który klika z nadmierną częstotliwością w reklamę nie generując przy tym konwersji.

4.1 Opis eksperymentu i danych użytych do analizy

Analiza została wykonana przy użyciu danych pochodzących z konkursu na najskuteczniejszy model klasyfikacyjny przeprowadzonego na platformie kaggle.com we współpracy z dostawcą usług reklamowych TalkingData, największej chińskiej platformy *big data* zajmującej się marketingiem internetowym²¹. Firma ta ogranicza problem *click fraudu* poprzez oflagowywanie adresów IP użytkowników, którzy generują nadmierną liczbę kliknięć i jednocześnie nie dokonują pożądaney konwersji. Poniższa analiza wykorzystuje wybrany podzbiór danych o kliknięciach w reklamę wraz z objaśnianą zmienną binarną, mówiącą o konwersji - pobraniu aplikacji mobilnej na telefon po danych kliknięciach. Do eksperymentu użyto dwóch zbiorów danych: treningowego i testowego o łącznej liczbie 1 032 340 obserwacji, których proces przygotowania został dokładnie opisany w rozdziale 5.1.

Głównym celem analizy jest zbudowanie modelu, który mógłby być wykorzystany przez dostawcę usług marketingu internetowego do klasyfikacji poszczególnych użytkowników na podstawie numerów identyfikacyjnych IP jako potencjalnych oszustów. Do budowy modelu zostały użyte metody uczenia maszynowego pozwalające na klasyfikację poszczególnych kliknięć jako takich, które mają odpowiednie prawdopodobieństwo wykonania konwersji po kliknięciu w reklamę. W następnej części zostaną opisane zmienne, które mają

²¹ TalkingData. (2018). TalkingData AdTracking Fraud Detection Challenge. Pobrane z: <https://www.kaggle.com/c/talkingdata-adtracking-fraud-detection/overview>. (Data dostępu: 08 maja 2021r.).

największy wpływ na wynik i pozwolą odpowiedzieć na pytanie, którzy użytkownicy generują kliknięcia, ale nie pobierają aplikacji.

4.2 Opis działania algorytmów użytych do predykcji

W pracy wykorzystane są dwa popularne modele uczenia maszynowego: Light Gradient Boosting Machine (LightGBM) oraz Extreme Gradient Boosting (XGBoost). Większa uwaga zostanie poświęcona pierwszemu modelowi (LightGBM), ze względu na większą dokładność wyników oraz znacznie szybszy czas trenowania samego modelu, co przekłada się na łatwość implementacji w praktyce biznesowej potencjalnego reklamodawcy i większą wartość dodaną w postaci ograniczenia strat. Model XGBoost został skonstruowany w celu porównania wyników i zweryfikowania większej efektywności modelu LightGBM. Algorytm LightGBM, oparty jest o wzmacniane gradientowo drzewa decyzyjne (Gradient Boosting Decision Trees) i został opracowany w ramach otwartej licencji przez Microsoft Research. Jest to stosunkowo nowy algorytm uczenia maszynowego, pierwsza stabilna wersja została opublikowana w roku 2017.²² W ostatnim czasie zyskał dużą popularność wśród społeczności ekspertów uczenia maszynowego z powodu wysokiej wydajności obliczeniowej i małego zużycia pamięci obliczeniowej przy dużych wolumenowo zbiorach danych posiadających dużą liczbę wymiarów. Kolejną jego zaletą, istotną przy niniejszej analizie, jest wysoka dokładność wyników²³. Ważnym aspektem działania tego algorytmu jest duża ilość danych potrzebnych do efektywnego uczenia, dlatego nie jest zalecane stosowanie go przy małych zbiorach danych, gdyż posiada wtedy skłonność do przetrenowywania. Dzięki swoim zaletom LightGBM zyskuje popularność względem XGBoost – drugiego algorytmu użytego w niniejszej analizie, również wykorzystującego technikę wzmacnianego gradientowego drzewa decyzyjnego.²⁴

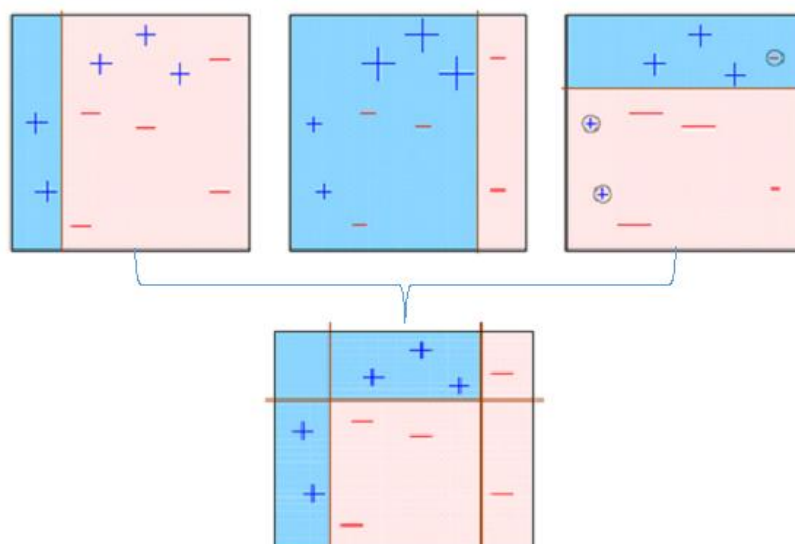
LightGBM, w ramach gradientowego wzmacniania drzew decyzyjnych, dodaje w sposób sekwencyjny (*ensemble learning*) kolejne predyktory, gdzie każdy następny poprawia dokładność klasyfikacji dopasowując się do błędu resztowego z poprzedniego predyktora. Dzięki tej procedurze sekwencyjnego dodawania kolejnych drzew decyzyjnych, będących słabymi klasyfikatorami, model minimalizuje błąd predykcji, osiągając dużą dokładność w klasyfikacji budując silny klasyfikator²⁵.

²² Microsoft Corporation. (2017). LightGBM documentation. Pobrane z: <https://github.com/microsoft/LightGBM>. (Data dostępu: 05 kwietnia 2021r.).

²³ Mandot, P. (2017). What is LightGBM, How to implement it? How to fine tune the parameters? Pobrane z: <https://medium.com/@pushkarmandot/https-medium-com-pushkarmandot-what-is-lightgbm-how-to-implement-it-how-to-fine-tune-the-parameters-60347819b7fc>. (Data dostępu: 09 maja 2021r.).

²⁴ Daoud (2019).

²⁵ Géron (2020, s. 212).

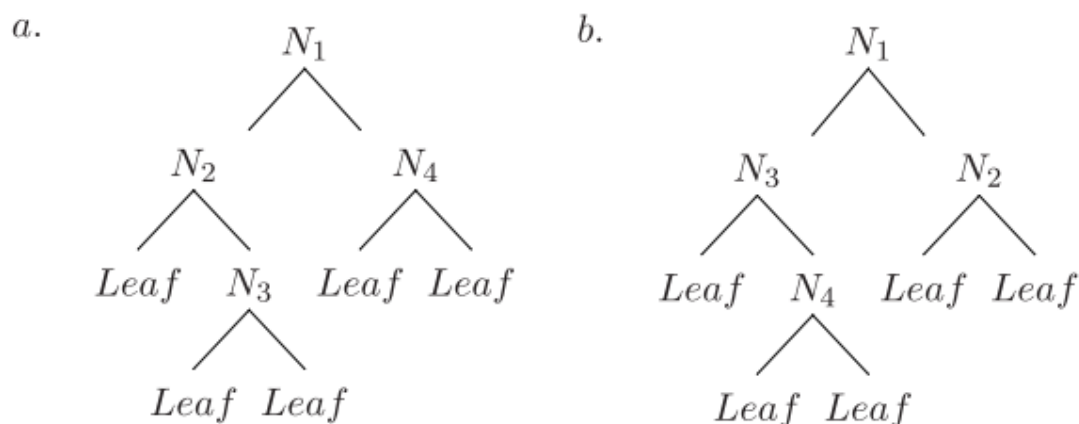


Rysunek 5: Tworzenie wzmacnianego klasyfikatora polegające na wytrenowaniu ostatecznego klasyfikatora przy pomocy wielu słabszych klasyfikatorów. Dzięki temu algorytmowi osiągane są dużo dokładniejsze wyniki. (Źródło: Analytics Vidhya, 2018. *Ensemble Learning | Ensemble Techniques*. Pobrane z: <https://www.analyticsvidhya.com/blog/2018/06/comprehensive-guide-for-ensemble-models/>. Data dostępu: 09 maja 2021r.)

Cechą wyróżniającą modelu LightGBM jest użycie nowatorskiej metody Gradient Based One-Side Sampling (GOSS), która wybiera próbki obserwacji zwiększających dokładność klasyfikacji na podstawie ich gradientów funkcji straty (*loss function*)²⁶. Algorytm w kolejnych sekwencjach uczenia zostawia próbki wyróżniające się dużym gradientem, czyli te które są niedopasowane, jednocześnie ograniczając obserwacje odznaczające się małym gradientem, czyli te dobrze dopasowane²⁷. Skutkiem zastosowania tego podejścia jest inna kolejność budowania drzew względem tego tradycyjnego, wykorzystywanego w XGBoost. W tradycyjnym drzewie decyzyjnym rozpoczynając od korzenia kolejne gałęzie są rozwijane w dół, rozwijając do końca daną gałąź. Natomiast w przypadku LightGBM rozpoczynając od korzenia każda kolejny liść jest wydzielany na podstawie najmniejszej wartości funkcji straty w całym dotychczas utworzonym drzewie. Podejście takie pozwala na szybszą budowę drzew o wysokiej skuteczności w klasyfikacji. Te różnice w kolejności budowanie kolejnych gałęzi drzewa decyzyjnego obrazuje Rysunek 6.

²⁶ Friedman (2001).

²⁷ Ke i in. (2017).



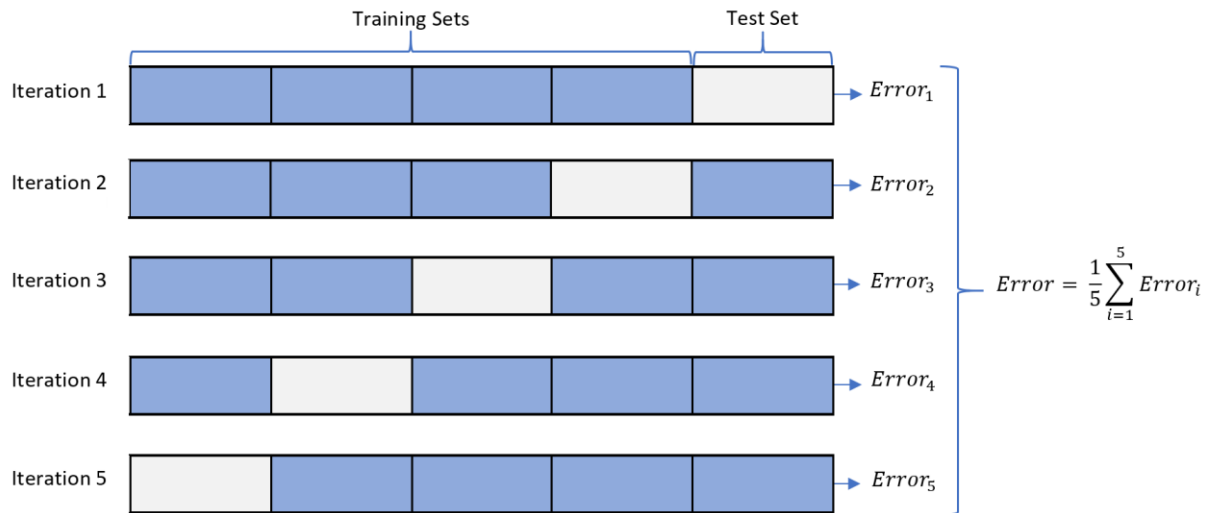
Rysunek 6: Schemat tworzenia drzewa tradycyjnego rozrastającego się w dół (a) i drzewa rozrastającego się według najlepszego dopasowania (b). Liczba w indeksie dolnym przy N mówi o kolejności tworzenia kolejnych rozgałęzień (Źródło: Shi, 2007, s. 4–5).

W przypadku budowy drzew bez ograniczeń, obydwa drzewa *a* i *b* miałyby taką samą budowę. Jednak w celu zwiększenia prędkości uczenia i zapobieganiu zjawiska przetrenowania stosuje się parametry ograniczające, m.in. maksymalną głębokość drzewa i minimalną ilość obserwacji w danym liściu, których skutkiem jest inna budowa kolejnych drzew decyzyjnych.

4.3 Wyznaczanie optymalnych hiperparametrów

Istotną zaletą algorytmu LightGBM jest ilość parametrów, które można dostosować w celu uzyskania jak największej kontroli nad procesem trenowania, w tym przypadku jest ich ponad 100. Niniejsza analiza została wykonana w środowisku Python 3.8.5 i w ramach niego został wykorzystany pakiet „lightgbm” na otwartej licencji firmy Microsoft dla modelu LightGBM i pakiet „xgboost” w przypadku algorytmu XGBoost. W celu uzyskania jak najlepszych rezultatów analiza została w dużym stopniu oparta na wynikach uzyskanych przez innych użytkowników platformy kaggle.com, którzy brali udział w konkursie i dzielili się swoimi sposobami podejścia do tego problemu klasyfikacyjnego. W implementacji obydwu algorytmów zostały dostosowane odpowiednie hiperparametry w celu lepszego dopasowania modelu i zmniejszenia skłonności do przetrenowywania się. Z powodu długiego czasu trenowania modeli użyto mechanizmu losowego doboru parametrów (RandomizedSearchCV z pakietu sci-kit learn) ze wskazanych przedziałów w każdym z optymalizowanych parametrów, która w oparciu o 5-krotny sprawdzian krzyżowy (ang. 5-fold Cross Validation) przeprowadza optymalizację wewnątrz zbioru treningowego. Sprawdzian krzyżowy wykorzystuje w tym przypadku uśredniony wynik z 5 iteracji trenowania modelu, każdej

wykonanej na innym podzbiore danych stanowiących w każdej iteracji 80% całego zbioru treningowego ($100\% * \left(1 - \frac{1}{5}\right) = 80\%$). Natomiast pozostałe 20% zbioru treningowego służy do testowania uzyskanych oszacowań.



Rysunek 7: Mechanizm działania 5-krotnego sprawdzianu krzyżowego (ang. 5-fold cross validation). Przy każdej z 5 iteracji walidowania modelu wykorzystywany jest inny, 20-procentowy podzbiór zbioru treningowego. Miara walidowanego modelu (w naszym przypadku AUC - pole pod krzywą ROC) jest uśrednionym wynikiem dla wszystkich 5 iteracji (Źródło: James, Witten, Hastie, Tibshirani, 2013, s. 184).

W kolejnym rozdziale zostaną przedstawione uzyskane wyniki modelowania uzyskane z użyciem optymalnych hiperparametrów wraz z oceną efektywności modeli.

5 Analiza efektywności klasyfikacji na przykładzie modelu LightGBM

W pierwszej części tego rozdziału zostaną opisane dane wykorzystane do analizy wraz z wskazaniem najważniejszych zmiennych służących do klasyfikacji i ich potencjalnym wpływem na model. W dalszej części przedstawione zostaną wyniki, ich interpretacja oraz ocena efektywności klasyfikacji. W końcowej części podsumowane zostaną wyniki i możliwości wykorzystania uzyskanego modelu w celu ograniczenia potencjalnych strat budżetów reklamowych na reklamę internetową reklamodawców.

5.1 Zbiór danych użyty do klasyfikacji

Z oryginalnych danych z serwisu kaggle.com²⁸ zawierających ponad 187 milionów kliknięć zostało wybranych 40 mln obserwacji, pierwszych w kolejności pod względem czasowym. Oryginalne dane pobrane w formacie CSV (Comma Separated Values) wyróżniały się bardzo niskim udziałem klasy pozytywnej, wynoszącym tylko 0,258% całego zbioru danych. Z tego względu została przeprowadzona procedura undersamplingu²⁹ dzięki której udział klasy pozytywnej wzrósł do 10%. Dane podzielono na dwa zbiory danych: treningowy i testowy w proporcji 82.7% do 17.3%. Wymienione dwa zbiory użyte do modelowania zostały wybrane zgodnie z przypisanym czasem kliknięcia (*click_time*). Okresy, w jakich zawierają się kliknięcia w zbiorach treningowym i testowym przedstawiają się następująco:

- zbiór treningowy – od 06.11.2017 godz. 14:32 do 07.11.2017 godz. 6:11 czasu UTC,
- zbiór testowy – od 07.11.2017 godz. 6:11 do 07.11.2017 godz. 9:39 czasu UTC.

Taki podział danych uwzględnia sekwencję działań i pozwala odwzorować sytuację, w której potencjalnie ten model miałby być używany do klasyfikacji fałszywych kliknięć.

Zbiór danych	Klasa negatywna	Klasa pozytywna	Łączna liczba obserwacji	Udział zbioru danych w całym zbiorze danych
Treningowy	768 258	85 362	853 620	82,7%
Testowy	160 848	17 872	178 720	17,3%
Łącznie	929 106	103 234	1 032 340	

Tabela 1 przedstawia liczbę obserwacji (kliknięć) w poszczególnych zbiorach użytych w procesie modelowania. Obydwa zbiory treningowy i testowy posiadają 10% udział klasy pozytywnej (Źródło: Opracowanie własne).

²⁸ TalkingData. (2018). TalkingData AdTracking Fraud Detection Challenge. Pobrane z: <https://www.kaggle.com/c/talkingdata-adtracking-fraud-detection/overview>. (Data dostępu: 08 maja 2021r.).

²⁹ Cateni, Colla, Vannucci (2014).

Spośród zmiennych zawartych w oryginalnym zbiorze danych do dalszych obliczeń wykorzystano następujące dane:

- *ip* - *zaszyfrowany* numer identyfikacyjny ip użytkownika, który dokonał kliknięcia w reklamę
- *app* – numer identyfikacyjny aplikacji, w którą kliknął użytkownik
- *os* – numer identyfikacyjny wersji oprogramowania urządzenia, z którego pochodzi kliknięcie w reklamę
- *device* – numer identyfikacyjny modelu urządzenia, z którego pochodzi kliknięcie
- *channel* – numer identyfikacyjny wydawcy reklamy
- *click_time* – czas (UTC) kliknięcia w reklamę
- *attributed_time* – czas (UTC) pobrania aplikacji pod warunkiem pobrania aplikacji przez użytkownika
- *is_attributed* – zmienna binarna, która jest objaśniana przez modele. Przyjmuje wartość 1 w przypadku pobrania aplikacji po kliknięciu w reklamę oraz 0 w przypadku kliknięcia, które nie skończyło się na jej pobraniu.

W celu przeprowadzenia analizy do oryginalnego zbioru danych, opisanego powyżej, dodano liczne zmienne określające pozostałą aktywność, które można podzielić na dziewięć głównych grup:

- Zmienne określające czas danego kliknięcia – *day*, *minute*, *hour* i *second*
- Grupowania według zmiennych *ip*, *app*, *device*, *os* oraz *channel*, obliczone dla tych grupowań: powtarzalność kliknięć w ciągu godziny, wariancja godzin oraz dni
- Średnia liczba kliknięć dla danej aplikacji przez użytkownika
- Częstość kliknięć dla danej aplikacji (zmienna *app*) i kanału (zmienna *channel*)
- Liczby unikalnych wartości *app*, *device*, *os* i *channel* dla danego numeru IP (zmienna *ip*)
- Skumulowana w czasie liczba kliknięć pochodzących od danego numeru IP przy tych samych wartościach *app*, *device*, *os* i *channel*
- Liczba kliknięć przy tych samych wartościach dla kombinacji zmiennych *app*, *device*, *os* i *channel* przed i po danym kliknięciu dla danego numeru IP
- Czas do następnego kliknięcia danego numeru IP (zmienna *ip*) przy danych zmiennych *app*, *device*, *os* oraz *channel*

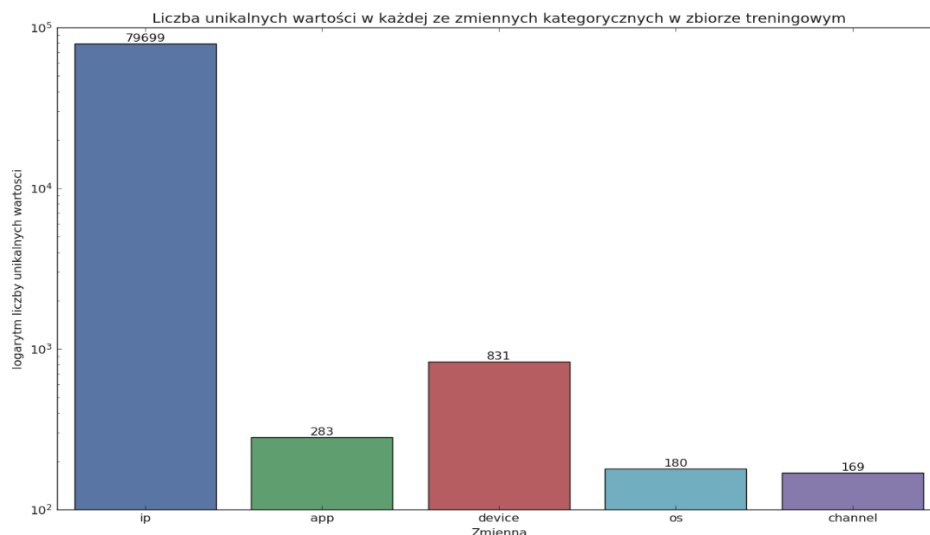
Powyższe zmienne zostały opracowane według autorskiego pomysłu oraz przy użyciu metod opisanych przez innego uczestnika konkursu na platformie kaggle.com³⁰. Te dodatkowe wymiary były obliczane wewnątrz każdego ze zbiorów (treningowym i testowym) w celu sprawdzenia efektywności modelu na kliknięciach, na których model klasyfikacyjny nie był jeszcze trenowany. Zmienne, które mają wyrażać powtarzalność kliknięć w danym zbiorze mają znacznie różniące się przedziały wartości ze względu na różnice w liczbie obserwacji pomiędzy zbiorami. W celu ujednolicenia tych przedziałów wartości wszystkie dodane wymiary zostały znormalizowane do przedziału [0,1]. Dodatkowo z powodu braku obsługi zmiennych kategoriycznych (dla zmiennych *channel*, *device* i *os*) przez algorytm XGBoost przygotowane zostały zmienne zero-jedynkowe przy pomocy metody OneHotEncoding z pakietu sci-kit learn, który generuje zmienne binarne przyjmujące wartość 1 w przypadku występowania danej kategorii dla danej obserwacji.

5.2 Wstępna analiza danych

Przed rozpoczęciem procesu modelowania, zbiór treningowy został poddany wstępnej eksploracji danych w celu znalezienia zależności, które mogłyby posłużyć lepszym wynikom klasyfikacyjnym modelu. W przypadku uczenia modeli maszynowych, tym bardziej tak skomplikowanych jak LightGBM i XGBoost istnieje dosyć spora trudność przy interpretowalności wyników oraz samodzielnym odnajdywaniu zależności. Te mogą być bardzo szczególne i niewidoczne przy wstępnej eksploracji danych. Z tego powodu niniejsza eksploracja danych jest stosunkowo krótka i przedstawia jedynie najważniejsze zależności zaobserwowane podczas eksploracji.

Rysunek 8 przedstawia ilości unikalnych wartości zmiennych kategoriycznych. W procesie budowy modelu zostały użyte wyłącznie trzy zmienne kategoriyczne *device*, *os* i *channel*. Zmienne *ip* i *app* zostały odrzucone z powodu oczekiwanego sposobu implementacji modelu, który zakłada jego wykorzystanie w celu rozpoznawania adresów *ip* użytkowników oraz dużej zmienności liczby unikalnych numerów IP i aplikacji w czasie. Ich użycie spowodowałoby pogorszającą się użyteczność modelu wraz z czasem.

³⁰ NanoMathias. (2018). Feature Engineering & Importance Testing. Pobrane z: <https://www.kaggle.com/nanomathias/feature-engineering-importance-testing#3.-Evaluating-Feature-Importance>. (Data dostępu: 10 maja 2021r.).



Rysunek 8: Liczba unikalnych wartości zmiennych kategorycznych *ip*, *app*, *device*, *os* i *channel* w zbiorze treningowym, przedstawiona na skali logarytmicznej, pokazuje dużą liczbę unikalnych numerów identyfikacyjnych (Źródło: Opracowanie własne przy użyciu pakietu *matplotlib* w środowisku *Python 3.8.5*).

Zmienne *channel* i *os* charakteryzują się dużymi rozbieżnościami w udziale klasy pozytywnej dla konkretnych kategorii, co obrazuje tabela 2.

<i>channel_id</i>	Liczebność	Udział klasy pozytywnej	<i>os_id</i>	Liczebność	Udział klasy pozytywnej
465	225	100%	61	928	100%
114	128	100%	29	3333	92%
408	111	100%	67	146	89%
274	7859	99%	24	7721	88%
419	954	98%	0	6526	85%
...
457	194	1%	53	5877	2%
417	1773	0%	46	1246	2%
125	5588	0%	42	1394	2%
364	2890	0%	23	7342	3%
420	84	0%	2	2515	1%

Tabela 2: Powyższe dwie tabele są ilustracją liczebności występowania i udziału klasy pozytywnej dla poszczególnych numerów identyfikacyjnych opisujących zmienne *channel* i *os*. W obydwu przypadkach wiersze w tabeli zostały wybrane spośród tych, których liczebność przekraczała 50. Dla każdej tabeli wybrano po 6 wierszy o najwyższym i najniższym udziale klasy pozytywnej w celach poglądowych – widoczna są duże rozbieżności w udziale klasy pozytywnej dla tych dwóch zmiennych (Źródło: Opracowanie własne).

Tak duże zróżnicowanie klasy pozytywnej i negatywnej dla tych dwóch zmiennych można interpretować następująco:

- Dla zmiennej *channel* oznaczającej poszczególnych wydawców reklamy, konkretni wydawcy odznaczają się mniejszą skutecznością reklamy bądź są w większym stopniu atakowani przez farmy botów, co powoduje mniejszy udział klasy pozytywnej.

- Dla zmiennej *os* oznaczającej konkretne wersje oprogramowania - farmy botów generujące nadmiarowe kliknięcia mogą wykorzystywać urządzenia o konkretnych wersjach oprogramowania, co przekłada się na mniejszy udział klasy pozytywnej.



Rysunek 9: Zmienna *app_AvgViewPerDistinct_ip* określa liczbę odwiedzin witryny danej aplikacji przez konkretny numer IP, na powyższym wykresie linia niebieska określa udział klasy pozytywnej zależnie od wartości opisanej zmiennej. Widać malejący udział klasy pozytywnej wraz z wzrostem liczby kliknięć w reklamę konkretnej aplikacji (Źródło: Opracowanie własne z wykorzystaniem pakietu *ggplot2* w środowisku R 4.0.2).

Powyższy wykres obrazuje zmniejszający się udział klasy pozytywnej wraz z wzrostem zmiennej *app_AvgViewPerDistinct_ip*, co pozwala podejrzewać, że użytkownicy identyfikujący się danymi numerami IP, częściej klikający w reklamę danej aplikacji, rzadziej dokonują pobrania aplikacji. Można znaleźć również podobne zależności m.in. dla zmiennych *ip_nextClick*, *ip_cumcount_app*, *ip_cumcount_os*.

5.3 Optymalizacja hiperparametrów modeli klasyfikacyjnych

W celu uzyskania jak najbardziej wiarygodnych wyników i wytrenowania najbardziej efektywnych modeli pod względem dokładności klasyfikacji dokonano optymalizacji hiperparametrów w przypadku obydwu modeli. Klasyfikacyjne modele uczenia maszynowego muszą posiadać odpowiednią umiejętność generalizacji reguł decyzyjnych przy jednoczesnej wysokiej dokładności klasyfikacji. Ta cecha zostaje zachowana dzięki doborowi optymalnych parametrów przy trenowaniu modeli. Zostały one wyznaczone według procedury opisanej w rozdziale 4.3. Postępując zgodnie z wcześniej opisanym schematem, uzyskano następujące optymalne rozwiązania dla modelu LightGBM:

- *learning_rate* = 0.1 – określa udział każdego drzewa w zespole. Przy niskich wartościach (mniejszych od 1) będzie wymagana większa liczba drzew, ale model najczęściej będzie lepiej generalizował wyniki³¹
- *num_leaves* = 45 – liczba liści w ostatecznym drzewie (domyślna wartość to 31)
- *max_depth* = 7 – maksymalna głębokość drzewa, ogranicza ryzyko przetrenowania modelu
- *colsample_bytree* \approx 0.616 – parametr używany do kontroli proporcji liczby kolumn dla losowej próbki
- *subsample* \approx 0.850 - współczynnik, który zwiększa szybkość obliczeń i ogranicza ryzyko przetrenowania modelu
- *min_data_in_leaf* = 1000 – minimalna liczba obserwacji zaklasyfikowanych do pojedynczego liścia ostatecznego drzewa decyzyjnego, pozwala ograniczyć ryzyko przetrenowania modelu³²

Odpowiedniki powyższych parametrów zostały ustawione dla modelu XGBoost na zbliżonych wartościach.

5.4 Struktura modelu

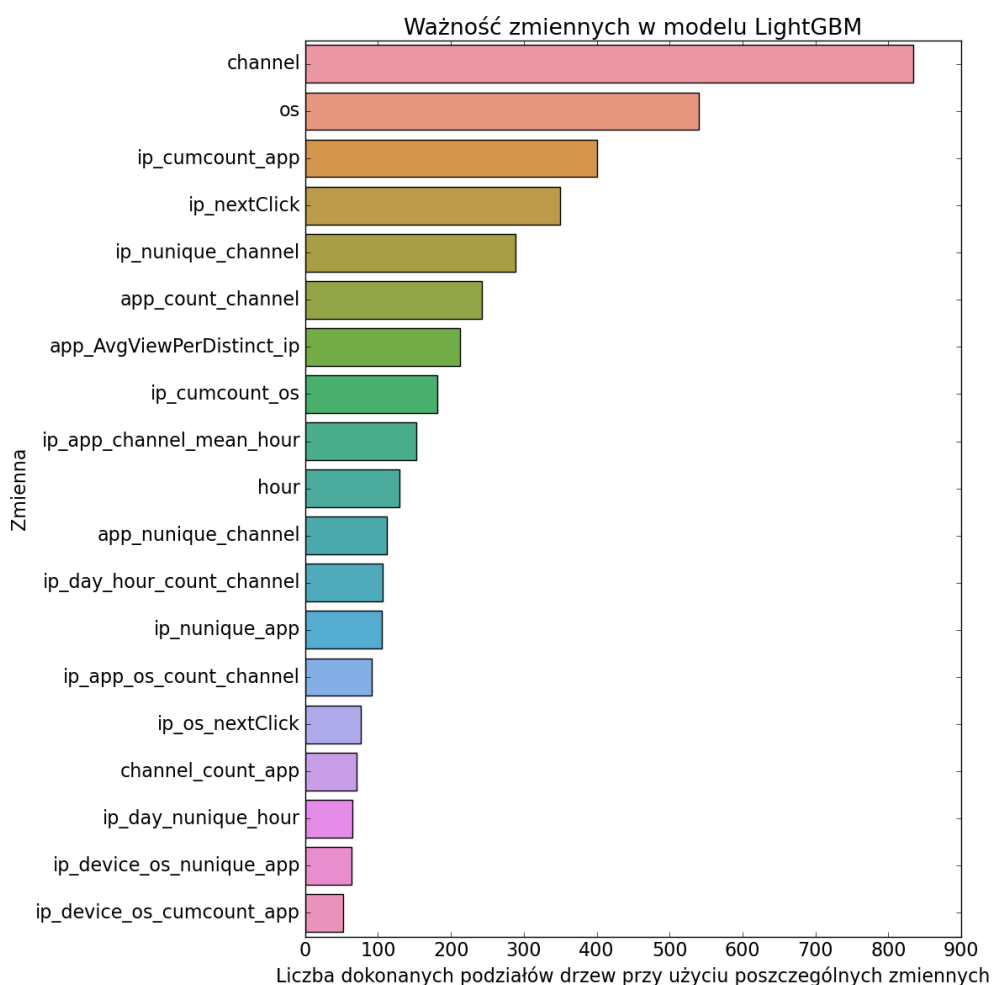
Po optymalizacji hiperparametrów uzyskano ostateczną postać modelu LightGBM, w której zmienne kategoryczne: *channel* i *os* wnoszą największą wartość do modelu. Zmienna *channel* określa numer identyfikacyjny wydawcy reklamy - konkretni wydawcy charakteryzują się zwiększonym współczynnikiem konwersji po kliknięciach. Charakterystyka zmiennej *os* jest podobna – użytkownicy konkretnych wersji oprogramowania urządzeń mobilnych częściej dokonują konwersji, może być to związane z używanymi wersjami oprogramowania przez tzw. farmy botów, które w sposób nadmierny klikają w reklamy nie generując przy tym pobrań aplikacji. Kolejne istotne dla modelu zmienne określają powtarzalność kliknięć pochodzących z danych numerów IP. Czwartą w kolejności jest zmienna *ip_nextClick*, która określa czas do następnego kliknięcia dla danego użytkownika.

Rysunek 10 przedstawia liczbę dokonanych podziałów przez poszczególne zmienne w zespole drzew tworzącym model. Wykres ten pokazuje duży wpływ zmiennych kategorycznych *channel* i *os*. Model jest zbudowany z 499 drzew decyzyjnych, które tworzą silny klasyfikator, dla którego zmienna *channel* stanowiła kryterium decyzyjne dla 834

³¹ Géron (2020, s. 213).

³² Microsoft Corporation. (2017). LightGBM documentation. Pobrane z: <https://github.com/microsoft/LightGBM>. (Data dostępu: 05 kwietnia 2021r.).

podziałów, natomiast zmienna *os* dla 540 podziałów. Według kryterium ilości podziałów w klasyfikatorze LightGBM w dalszej kolejności silne reguły klasyfikacyjne tworzą zmienne opisujące powtarzalność kliknięć w reklamę i odstęp czasowy pomiędzy kolejnymi kliknięciami (*ip_nextClick*) oraz skumulowane w czasie liczby kliknięć dla danych kombinacji numerów IP, aplikacji i kanału reklamowego.

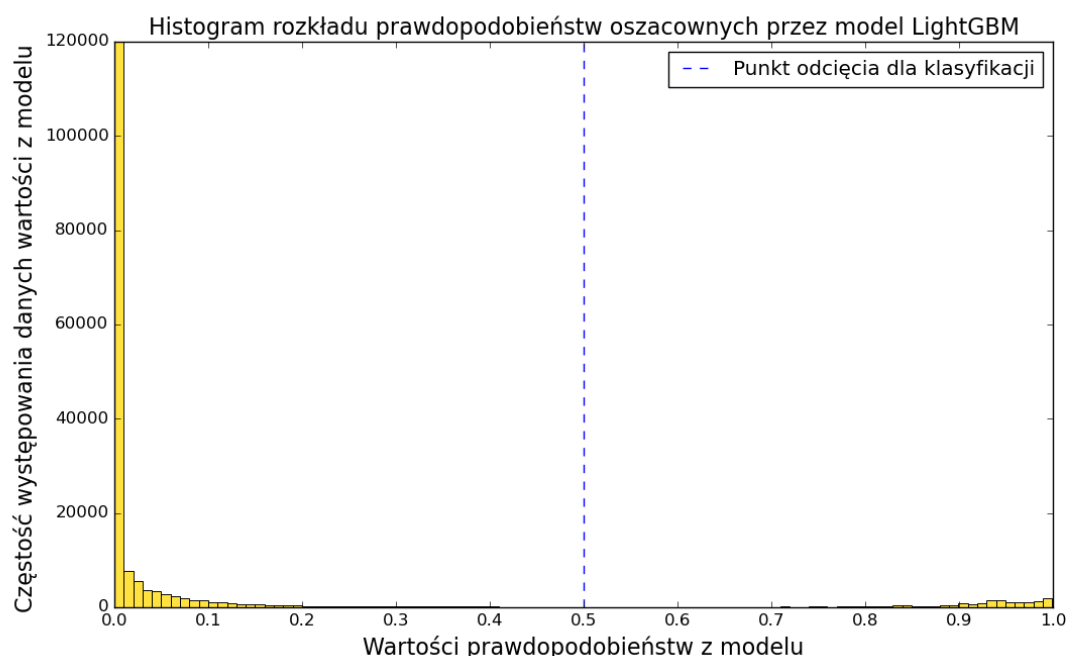


Rysunek 10: Liczba dokonanych podziałów drzew w modelu LightGBM przez zmienne użyte w modelu. Ze względu na dużą liczbę zmiennych w modelu rysunek przedstawia zmienne, które dokonały co najmniej 50 takich podziałów. Widać znaczący wpływ zmiennych kategoriycznych *channel* i *os* na klasyfikację (Źródło: Opracowanie własne z wykorzystaniem pakietu *matplotlib* w środowisku *Python 3.8.5*).

5.5 Wyniki klasyfikacji

Wyniki klasyfikacji dla obydwu modeli zostały oparte o predykcje na podstawie zbioru testowego zawierającego 178 720 obserwacji, wśród których około 10% stanowiła klasa pozytywna. Obydwa modele prognozują prawdopodobieństwo przyporządkowania do klasy pozytywnej, na którego podstawie jest dokonywana klasyfikacja binarna. Przy pomocy tych

prawdopodobieństw była dokonywana decyzja polegająca na przypisaniu do klasy negatywnej lub pozytywnej zależnie od przekroczenia wartości punktu odcięcia, która wynosiła 0.5. Rysunek 11 przedstawia rozkład prawdopodobieństw oszacowanych przez model LightGBM, w na podstawie których jest dokonywana klasyfikacja.



Rysunek 11: Rozkład wyników z modelu LightGBM. Powyższy wykres reprezentuje rozkład wyników prawdopodobieństw otrzymanych z modelu LightGBM, przerywaną linią został oznaczony punkt odcięcia, który służył za regułę decyzyjną. Obserwacje, które przekroczyły ten punkt odcięcia, równy 0.5 zostały zaklasyfikowane do grupy pozytywnej, oznaczającej obecność konwersji po kliknięciu w reklamę (Źródło: Opracowanie własne z wykorzystaniem pakietu matplotlib w środowisku Python 3.8.5).

Powyższy rysunek wskazuje na dużą asymetrię rozkładu prawdopodobieństw prognozowanych przez model – 90.6% obserwacji posiada wartość mniejszą niż 0.5, co jest spowodowane niezbalansowanym zbiorem testowym. Wartość 0.5 przekracza 9.4% obserwacji, spośród których 85.6% jest większych niż 0.8, co może wskazywać na wysokie umiejętności rozpoznawania klasy pozytywnej w zbiorze testowym.

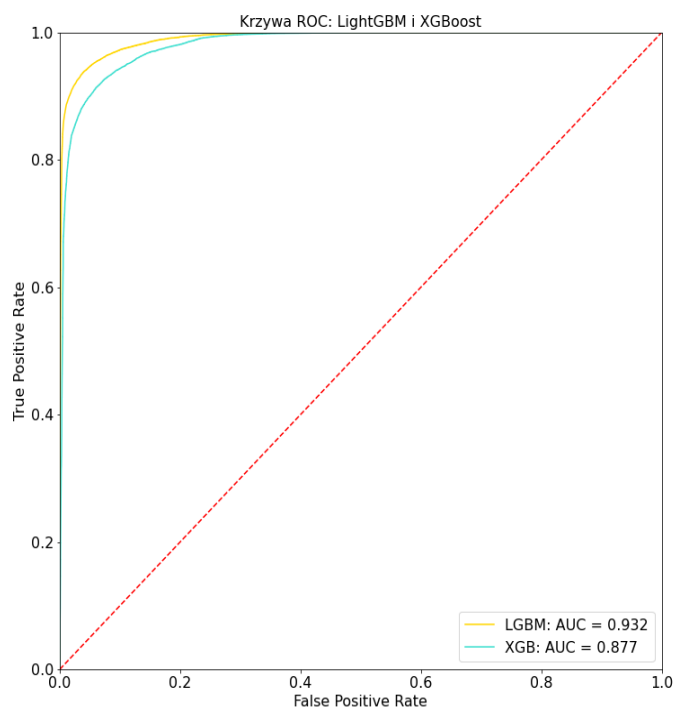
5.6 Walidacja modelu

5.6.1 Krzywa ROC i współczynnik AUC

Przy walidacji modeli i porównywaniu dokładności klasyfikacji pomiędzy wytrenowanymi modelami używano najczęściej wykorzystywanej metryki w przypadku problemów klasyfikacji binarnej, czyli obszaru pod krzywą charakterystyki roboczej odbiornika (*receiver operating curve*) dalej określanego jako AUC (*area under the curve*).³³

³³ James, Witten, Hastie, Tibshirani (2013, s. 148).

Wykresy krzywych ROC dla obydwu modeli przedstawia Rysunek 12. W prawym dolnym rogu rysunku zostały przedstawione wartości wskaźnika AUC dla tych krzywych. Wysoki wynik AUC dla modelu LightGBM równy 0.932 pokazuje dużo większą skuteczność względem modelu XGBoost (AUC = 0.877).



Rysunek 12: Krzywa ROC dla modeli LightGBM i XGBoost - pokazuje relację między odsetkiem prawdziwie pozytywnych klasyfikacji (pełnością) modelu i odsetkiem fałszywie pozytywnych (Źródło: Opracowanie własne przy pomocy pakietu matplotlib w środowisku Python 3.8.5)

5.6.2 Analiza macierzy pomyłek

W celu określenia dokładności i specyfiki klasyfikacji na obydwu modelach zostały użyte macierze pomyłek zbudowane według poniższego schematu z Tabela 3.

Macierz pomyłek	Klasyfikacja negatywna	Klasyfikacja pozytywna
Rzeczywisty stan negatywny	Obserwacje prawdziwie negatywne (PN)	Obserwacje fałszywie pozytywne (FP)
Rzeczywisty stan pozytywny	Obserwacje fałszywie negatywne (FN)	Obserwacje prawdziwie pozytywne (PP)

Tabela 3: Użyty schemat macierzy pomyłek (Źródło: Géron, 2020, s. 110).

Dla modelu LightGBM otrzymana macierz pomyłek ma postać:

Macierz pomyłek: LightGBM	Klasyfikacja negatywna	Klasyfikacja pozytywna
Rzeczywisty stan negatywny	159 570	1 278
Rzeczywisty stan pozytywny	2 296	15 576

Tabela 4: Macierz pomyłek dla modelu LightGBM obliczona na podstawie zbioru testowego (Źródło: Opracowanie własne w środowisku Python 3.8.5).

Na podstawie powyższej macierzy dla modelu LightGBM zostały obliczone następujące charakterystyki efektywności klasyfikacji modelu:

1. $Dokładność = \frac{(PP+PN)}{\text{liczba obserwacji}} = 0.98$ – uzyskano bardzo wysoki wynik dokładności mówiący o udziale obserwacji poprawnie sklasyfikowanych przez model. Jednak ze względu na niezbalansowanie zbioru testowego wynik ten nie ocenia wiarygodnie efektywności modelu.
2. $Precyzja = \frac{PP}{PP+FP} = 0.92$ - dodatnia wartość predykcyjna
3. $Czułość = \frac{PP}{PP+FN} = 0.87$ - odsetek prawdziwie pozytywnych obserwacji
4. $Swoistość = \frac{PN}{PN+FP} = 0.99$ - odsetek prawdziwie negatywnych obserwacji³⁴

Ważność przedstawionych powyżej miar wydajności modeli zależy od celu wykorzystania modelu. Niniejsza praca poświęcona jest problemowi fałszywych kliknięć i ich wykrywaniu, dlatego opisane modele powinny mieć jak największą skuteczność w przypadku klasy negatywnej. Dlatego istotną miarą w przypadku tego modelu jest miara swoistości, utrzymana na bardzo wysokim poziomie 99% w przypadku obydwu modeli, co pozwala na stwierdzenie o dobrej jakości modeli. Należy jednak zauważyć, że miary trafności klasyfikacji klasy pozytywnej (precyzja i czułość) również zachowują wysokie wartości oscylujące wokół wartości 0.9 dla głównego modelu LightGBM.

Model	Dokładność	Precyzja	Czułość	Swoistość
LightGBM	0.98	0.92	0.87	0.99
XGBoost	0.97	0.88	0.76	0.99

Tabela 5: Porównanie miar wydajności klasyfikacji dla modeli LightGBM i XGBoost (Źródło: Opracowanie własne w środowisku Python 3.8.5).

³⁴ James i in. (2013, s. 146).

Powyżej przedstawione wyniki pokazują wyższą efektywność modelu LightGBM względem XGBoost w przypadku skuteczności klasyfikacji dla klasy pozytywnej. Jednocześnie swoistość dla obydwu modeli na poziomie 0.99 pokazuje wysokie umiejętności predykcji klasy negatywnej dla obydwu modeli.

5.7 Przykładowe zastosowanie modelu w praktyce biznesowej

Firmy, które zajmują się problemem wykrywania fałszywych kliknięć ograniczają straty reklamodawców poprzez oflagowywanie użytkowników generujących dużą liczbą kliknięć przy jednoczesnej niskiej skłonności do wykonywania konwersji. Takie podejście również miała zamiar zastosować platforma TalkingData, która udostępniła dane o kliknięciach w ramach konkursu na platformie kaggle.com³⁵. Mechanizm ten może działać z dużą skutecznością, dzięki zastosowaniu reguł decyzyjnych opracowanych w oparciu o model klasyfikacyjny działający według algorytmów uczenia maszynowego. Na podstawie powyżej opisanego modelu klasyfikacyjnego LightGBM jest możliwe wypracowanie reguły decyzyjnej, mającej na celu oznaczanie użytkowników identyfikujących się numerami adresów, którzy generują fałszywe kliknięcia i doprowadzają swoją nieuczciwą działalnością do strat budżetowych reklamodawców.



Rysunek 13: Powyższy wykres rozproszenia średnich wartości oszacowanych prawdopodobieństw wykonania konwersji i liczby kliknięć dla danego numeru IP jest wizualizacją reguły decyzyjnej, która selekcjonuje adresy IP o wysokiej liczbie kliknięć przy niskim średnim prawdopodobieństwie wykonania konwersji (mniejszym niż 0.1) i wysokiej liczbie kliknięć (większej niż 50) dla zbioru testowego obejmującego obserwacje z przedziału czasowego o długości 2h 26 min. (Źródło: Opracowanie własne przy pomocy pakietu ggplot2 w środowisku R 4.0.2).

³⁵ TalkingData. (2018). TalkingData AdTracking Fraud Detection Challenge. Pobrane z: <https://www.kaggle.com/c/talkingdata-adtracking-fraud-detection/overview>. (Data dostępu: 08 maja 2021r.).

Rysunek 13 obrazuje zastosowanie przykładowej reguły decyzyjnej, która wybiera adresy IP podejrzane o generowanie fałszywych kliknięć, charakteryzujące się:

- niskim średnim prawdopodobieństwem wykonania konwersji oszacowanym przez model wynoszącym mniej niż 0.1
- wysoką liczbą kliknięć w reklamę wynoszącą więcej niż 50.

Na podstawie tej prostej reguły następuje klasyfikacja konkretnych numerów IP:

- prawidłowych - zachowujących się prawidłowo
- podejrzanych - użytkowników, którzy prawdopodobnie generują fałszywe kliknięcia.

W ten sposób, spośród 45 613 unikalnych adresów IP w zbiorze testowym, 93 użytkowników zostało zaklasyfikowanych jako podejrzanych o generowanie fałszywych kliknięć. Zastosowanie takiej reguły decyzyjnej przy klasyfikacji konkretnych adresów IP jako oszustów pozwala na oznaczenie 6.56% kliknięć w zbiorze testowym jako fałszywych. Reklamodawca lub wydawca reklamy internetowej, korzystający z powyższej reguły decyzyjnej może znacznie ograniczyć ruch internetowy pochodzący z nieuczciwego źródła, oraz negatywne skutki oszustw reklamowych w postaci utraty budżetu reklamowego, jednocześnie polepszając wiarygodność statystyk używanych do mierzenia kampanii.

6 Wnioski

W pracy przedstawiono wykorzystanie algorytmów uczenia maszynowego w celu predykcji dokonania konwersji na podstawie zmiennych opisujących kliknięcia w reklamę wykonane na stronie internetowej. Uzyskane wyniki z modelu wskazują na wysoką skuteczność klasyfikacji, co pokazuje wysoka wartość współczynnika AUC na poziomie 0.932. Na podstawie zbudowanego modelu została przedstawiona przykładowa reguła decyzyjna wspomagająca proces decyzyjny, mająca na celu oznaczenie oszustów reklamowych generujących nadmierną liczbę kliknięć w czasie przy niskiej liczbie konwersji. Przedstawiona metoda oflagowywania oszustów reklamowych jest bardzo uproszczona i posiada duże możliwości rozwoju. Opisana reguła decyzyjna stanowi jednak dużą wartość dla korzystającego z niej reklamodawcy. Dzięki odfiltrowywaniu ruchu na stronie z potencjalnych oszustów reklamowych pozwala na ograniczenie negatywnych skutków zjawiska fałszywych kliknięć w postaci strat budżetu reklamowego i fałszowania wyników kampanii.

W celu opracowania bardziej efektywnej reguły decyzyjnej wskazane jest zbudowanie modelu prognozującego dokonanie konwersji po kliknięciu w reklamę na podstawie pełnego zbioru danych. Jednak to wymaga bardzo dużej mocy obliczeniowej, co może być problemem w przypadku podmiotów chcących ograniczyć straty, ale nie posiadających odpowiednio dużego budżetu na implementację takiego rozwiązania. Dodatkową możliwością rozwoju jest udoskonalenie reguły decyzyjnej poprzez dodanie kolejnych kryteriów oceny, mogących poprawić skuteczność klasyfikacji potencjalnych oszustów reklamowych.

Bibliografia

Publikacje naukowe

- Breiman, L. (1997). Arcing the Edge. *Technical Report*. (496).
- Cateni, S., Colla, V., Vannucci, M. (2014). A method for resampling imbalanced datasets in binary classification tasks for real-world problems. *Neurocomputing*, 135, 32–41. DOI: <https://doi.org/10.1016/j.neucom.2013.05.059>.
- Daoud, E. A. (2019). Comparison between XGBoost, LightGBM and CatBoost Using a Home Credit Dataset. *International Journal of Computer and Information Engineering*. (13 (1)), 6–10. DOI: <https://doi.org/10.1287/mksc.1080.0397>.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29(5), 1189–1232. DOI: <https://doi.org/10.1214/aos/1013203451>.
- Géron, A. (2020). *Uczenie maszynowe z użyciem Scikit-Learn i TensorFlow* (wyd. 2). Gliwice: Helion.
- James, G., Witten, D., Hastie, T., Tibshirani, R. (red.). (2013). *An Introduction to Statistical Learning: with Applications in R* (wyd. 7). New York: Springer.
- Jothi, G. (2019). A study on advantages and disadvantages of traditional marketing and digital marketing. *International Journal of Research and Analytical Reviews*.
- Karjaluoto, H., Mustonen, N., Ulkuniemi, P. (2015). The role of digital channels in industrial marketing communications. *Journal of Business & Industrial Marketing*, 30(6), 703–710. DOI: <https://doi.org/10.1108/JBIM-04-2013-0092>.
- Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., . . . Liu, T.-Y. (2017). LightGBM: A Highly Efficient Gradient Boosting Decision Tree. W: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (red.), *Advances in Neural Information Processing Systems* (nr. 30). Curran Associates, Inc.
- Nair, K., Gupta, R. (2021). Application of AI technology in modern digital marketing environment. *World Journal of Entrepreneurship, Management and Sustainable Development, ahead-of-print*(ahead-of-print). DOI: <https://doi.org/10.1108/WJEMSD-08-2020-0099>.
- Oentaryo, R., Lim, E.-P., Finegold, M., Lo, D., Zhu, F., Phua, C., . . . Berrar, D. (2014). Detecting click fraud in online advertising: a data mining approach. *Journal of Machine Learning Research*. (15), 99–140.
- Shi, H. (2007). *Best-first Decision Tree Learning* (Master of Science). University of Waikato, Hamilton, New Zealand.
- Wilbur, K., Zhu, Y. (2008). Click fraud. *Marketing Science*. (28), 293–308. DOI: <https://doi.org/10.1287/mksc.1080.0397>.

Raporty badawcze

- Cramer-Flood, E. (2020). Global digital ad spending update Q2 2020. Pobrane z: baza publikacji eMarketer pod adresem <https://content-na1.emarketer.com/global-digital-ad-spending-update-q2-2020>. (Data dostępu: 4 kwietnia 2021r.)

Cramer-Flood, E. (2020b). Global Ecommerce 2020: Ecommerce Decelerates amid Global Retail Contraction but Remains a Bright Spot. Pobrane z: baza publikacji eMarketer pod adresem <https://content-na1.emarketer.com/global-ecommerce-2020>. (Data dostępu: 4 kwietnia 2021r.)

Perrin, N. (2020). Digital Ad Fraud 2020. Pobrane z: baza publikacji eMarketer pod adresem <https://content-na1.emarketer.com/digital-ad-fraud-2020>. (Data dostępu: 4 kwietnia 2021r.)

Spis rysunków

Rysunek 1: Globalny udział wydatków na reklamę cyfrową.....	8
Rysunek 2: Globalny udział handlu elektronicznego i jego roczny wzrost	10
Rysunek 3: Najwięksi dostawcy usług reklamy cyfrowej.....	13
Rysunek 4: Procentowy udział źródeł fałszywych kliknięć	17
Rysunek 5: Schemat tworzenia wzmacnianego klasyfikatora.....	21
Rysunek 6: Porównanie sposobów budowy drzewa decyzyjnego.	22
Rysunek 7: Mechanizm działania 5-krotnego sprawdzianu krzyżowego.	23
Rysunek 8: Liczba unikalnych wartości zmiennych kategorycznych	27
Rysunek 9: Udział klasy pozytywnej przy zmianie zmiennej <i>app_AvgViewPerDistinct_ip</i> . ..	28
Rysunek 10: Liczba dokonanych podziałów drzew w modelu LightGBM.....	30
Rysunek 11: Rozkład wyników z modelu LightGBM.	31
Rysunek 12: Krzywe ROC dla modeli LightGBM i XGBoost	32
Rysunek 13: Wykres rozproszenia średnich wartości oszacowanych prawdopodobieństw wykonania konwersji i liczby kliknięć dla danego numeru IP	34

Spis tabel

Tabela 1: Liczba obserwacji w poszczególnych zbiorach użytych w procesie modelowania.	24
Tabela 2: Liczebność i udział klasy pozytywnej dla poszczególnych numerów identyfikacyjnych opisujących zmienne <i>channel</i> i <i>os</i>	27
Tabela 3: Schemat macierzy pomyłek	32
Tabela 4: Macierz pomyłek dla modelu LightGBM.....	33
Tabela 5: Porównanie miar wydajności klasyfikacji dla modeli LightGBM i XGBoost	33

Streszczenie

Zjawisko fałszywych kliknięć stanowi poważny problem w marketingu cyfrowym, powodując znaczące straty budżetowe dla reklamodawców i przyczyniając się do niskiej wiarygodności miar efektów kampanii. Celem niniejszej pracy było skonstruowanie narzędzia opartego o metody uczenia maszynowego, umożliwiającego ograniczenie tego zjawiska, poprzez identyfikację potencjalnych oszustów reklamowych. Do budowy narzędzia zostały wykorzystane dane pochodzące od chińskiego dostawcy usług marketingowych, dotyczące kliknięć w reklamę i pobrania aplikacji po danych kliknięciach, określanych jako konwersję. Na potrzeby analizy przekształcono oryginalny zbiór danych w taki sposób, by zawierał informacje o powtarzalności kliknięć w czasie. Na ich podstawie skonstruowany został model klasyfikujący, czy zostanie wykonana konwersja po danych kliknięciach. W tym celu wykorzystano algorytm LightGBM oraz dla porównania skuteczności tego modelu użyto także metody XGBoost. Wyniki z modelu klasyfikacyjnego w postaci średnich wartości prawdopodobieństw wykonania konwersji po kliknięciu w reklamę dla użytkownika oraz ilość kliknięć wygenerowanych przez użytkownika w czasie, posłużyły jako dwa kryteria decyzyjne tego narzędzia. W ten sposób zbudowana reguła decyzyjna pozwoliła na oflagowanie użytkowników jako potencjalnych oszustów, przyczyniając się do oznaczenia znacznej części kliknięć jako fałszywych. Wykorzystanie tej reguły decyzyjnej do odfiltrowywania oszustów reklamowych w dłuższym okresie pozwala na ograniczenie negatywnych skutków zjawiska fałszywych kliknięć.