

Analiza potencjału rynku indyjskiego dla smartfona OnePlus 11R 5G

ULADZISLAU ZHERABIATSYEU^{1,*}

¹ Politechnika Wrocławska, Wydział Informatyki i Telekomunikacji, Informatyka Stosowana, MSiD Lab czwartek 17:05 TN

* 266959@student.pwr.edu.pl

Compiled June 12, 2023

W raporcie jest analizowany potencjał indyjskiego rynku dla smartfona OnePlus 11R 5G poprzez analizę danych. Badanie obejmuje eksperyment, który bada recenzje Amazon, wyniki analizy sentymentu, macierz korelacji i inne wskaźniki, takie jak liczba ocen według koloru, rozkład rozmiarów pamięci RAM, rozkład kolorów i najpopularniejsze słowa według sentymentu. Badanie wykazało, że średnia ocena dla obu kolorów jest prawie taka sama, z wysokim wyborem Galactic Silver w marcu 2023 roku. Ponadto badanie podkreśla znaczenie funkcji, takich jak pamięć RAM, pamięć masowa i zweryfikowany status, oraz ich korelację z kategoriami sentymentu. Wyniki sugerują, że smartfon otrzymał ogólnie pozytywne opinie, ale istnieją pewne negatywne recenzje dotyczące kwestii związanych z telefonem, problemów i doświadczeń związanych z zakupem. Wreszcie, badanie ujawnia, że poziom oceny 5.0 początkowo wzrósł, ale następnie spadł, a nastroje radości i neutralności są bardziej powszechne wśród zweryfikowanych recenzji.

1. WSTĘP

Problemem poddawanym badaniom jest zwiększenie atrakcyjności produktu "OnePlus 11R 5G" poprzez zidentyfikowanie kluczowych czynników wpływających na preferencje i decyzje zakupowe konsumentów. Celem raportu jest przeprowadzenie analizy danych dotyczących opinii klientów na temat telefonu oraz określenie, które cechy produktu są najważniejsze dla nabywców. W szczególności, analiza skupi się na takich parametrach jak pojemność pamięci, ilość RAM czy status zweryfikowanego nabywcy. Raport ma pomóc producentowi w podejmowaniu trafnych decyzji biznesowych związanych z ulepszaniem i promocją produktu. Jest adresowany do osób odpowiedzialnych za rozwój produktu oraz strategię marketingowe w firmie OnePlus.

2. ZBIÓR DANYCH I JEGO PRZETWARZANIE

A. Zbiór danych

Dane opinii o produkcie zostały zebrane ze strony internetowej amazon.in, dotyczącej smartfona OnePlus 11R 5G. Zawartość zbioru danych obejmuje szereg informacji na temat opinii użytkowników, takich jak tytuł, treść, ocena, data wystawienia opinii, ilość RAM-u, ilość pamięci, kolor oraz czy opinia jest zweryfikowana.

B. Przetwarzanie wstępne

Konstrukcja modelu wymagała odpowiedniego przetworzenia danych wejściowych:

1. Usunięcie wierszy zawierających wartości puste.

2. Sprawdzenie typów danych w poszczególnych kolumnach w celu potwierdzenia możliwości ich konwersji na liczby.
3. Konwersja kolumn z treścią i tytułem na typ string.
4. Konwersja kolumny z datą na typ daty, indeksacja, sortowanie całej modeli względem daty.
5. Usunięcie duplikatów.

3. ANALIZA EKSPLORACYJNA

A. Kolor

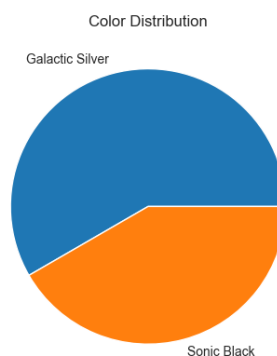


Fig. 1. Ilość opinii względem koloru

Color	Galactic Silver	Sonic Black
Ilość	629	451

Table 1. Ilość opinii względem koloru

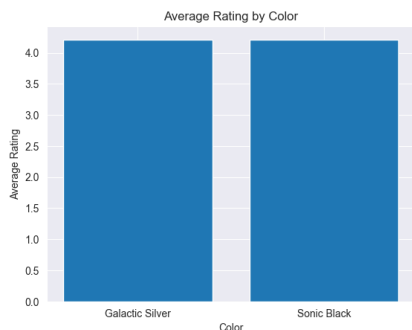


Fig. 2. Średnia Ocena względem koloru

Color	Galactic Silver	Sonic Black
Średnia ocena	4.20	4.21

Table 2. Średnia Ocena względem koloru

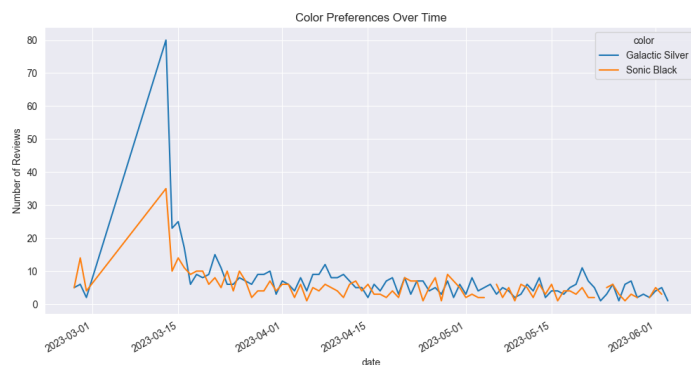


Fig. 3. Popularność kolorów względem czasu

Na podstawie statystyk można wnioskować, że klienci wykazują zrównoważone preferencje co do kolorów produktu. W analizowanym okresie najwięcej sprzedaży przypada na kolory Galactic Silver oraz Sonic Black, przy czym średnia ocena użytkowników dla obu tych kolorów jest bardzo zbliżona i wynosi odpowiednio 4.20 i 4.21.

Warto również zauważyć, że w dniu 14-03-2023 r. nastąpił duży wzrost popularności koloru Galactic Silver, który został najczęściej wybierany przez klientów. W kolejnych dniach preferencje kolorystyczne były bardziej zróżnicowane, a wybór między kolorem Sonic Black a Galactic Silver był coraz bliższy.

B. RAM

RAM	8GB	16GB
Ilość	600	447

Table 3. Ilość opinii względem ilości RAM

Distribution of RAM Sizes

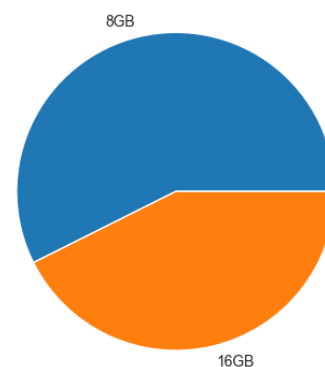


Fig. 4. Ilość opinii względem ilości RAM

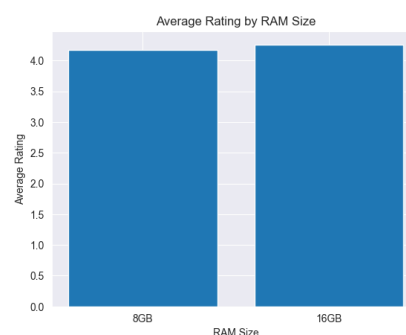


Fig. 5. Średnia Ocena względem ilości RAM

RAM	8GB	16GB
Średnia ocena	4.17	4.26

Table 4. Średnia Ocena względem ilości RAM

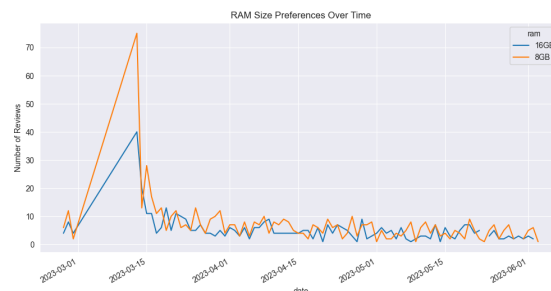


Fig. 6. Popularność RAM względem czasu

Na podstawie statystyk można wywnioskować, że zarówno 8GB, jak i 16GB pamięci RAM są popularne wśród konsumentów, ale średnia ocena urządzeń z 16GB RAM jest nieco wyższa. Dodatkowo, wydaje się, że na rynku występuje pewna niepewność przy wyborze modelu telefonu na podstawie wielkości RAM-u, wielu konsumentów początkowo decyduje się na wersję 8GB, ale ostatecznie dokonuje wyboru na podstawie innych czynników.

C. Zweryfikowany status

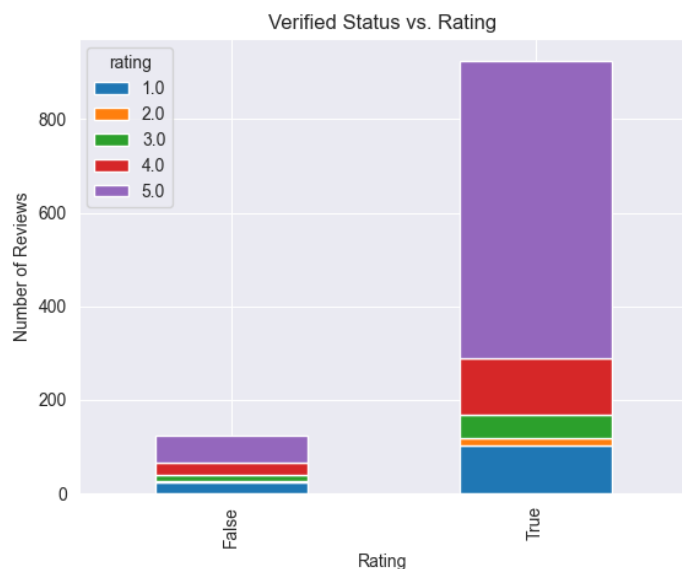


Fig. 7. Stosunek zweryfikowanego statusu do oceny

	Rating						
Verified	1.0	2.0	3.0	4.0	5.0	Total	% Verified
False	24	1	15	25	57	122	19.67%
True	102	15	51	120	637	925	80.33%
Total	126	16	66	145	694	1047	
% False	19.67%	0.82%	12.30%	20.49%	46.72%		
% True	11.03%	1.62%	5.51%	12.97%	68.86%		

Table 5. Stosunek zweryfikowanego statusu do oceny

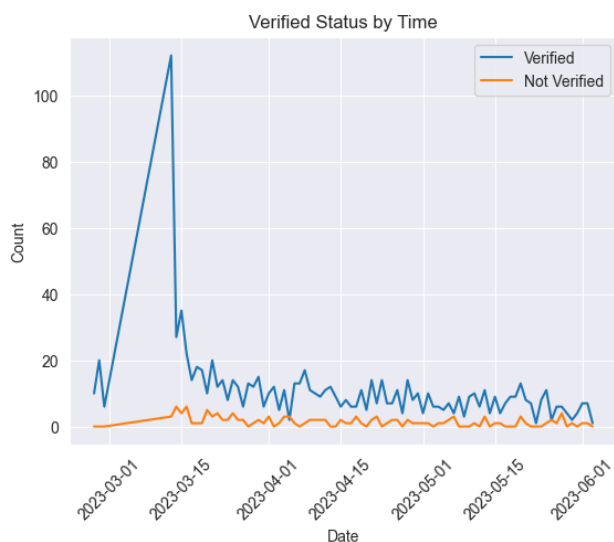


Fig. 8. Ilość opinii opinii względem statusu weryfikacji i czasu

Na podstawie statystyk można wyciągnąć kilka wniosków odnośnie produktu i jego oceny. Przede wszystkim, widoczny jest wyraźny wzrost liczby zweryfikowanych opinii w porównaniu do niezwyfikowanych w dniu 13.03.2023. Jest to pozytywny sygnał dla produktu, ponieważ większa ilość zweryfikowanych opinii sugeruje większą wiarygodność i zaufanie konsumentów.

Patrząc na stosunek zweryfikowanego statusu do oceny, można zauważyć, że większość opinii (80,33 procentów) jest zweryfikowana, co również wpływa na wiarygodność produktu. Ponadto, najwięcej opinii dotyczy oceny 5.0, co może oznaczać, że produkt cieszy się dużym uznaniem wśród klientów.

Jeśli chodzi o procentowy stosunek zweryfikowanego statusu do oceny, to warto zwrócić uwagę na fakt, że najwyższy procent opinii zweryfikowanych występuje przy ocenie 5.0 (68,86 procentów), co dodatkowo podkreśla pozytywny wizerunek produktu. Z kolei najwyższy procent opinii niezwyfikowanych dotyczy oceny 5.0 (46,72 procentów), co może budzić pewne wątpliwości co do ich autentyczności.

D. Popularność i satysfakcja

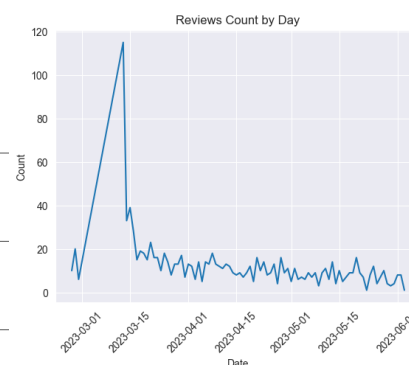


Fig. 9. Ilość opinii względem czasu

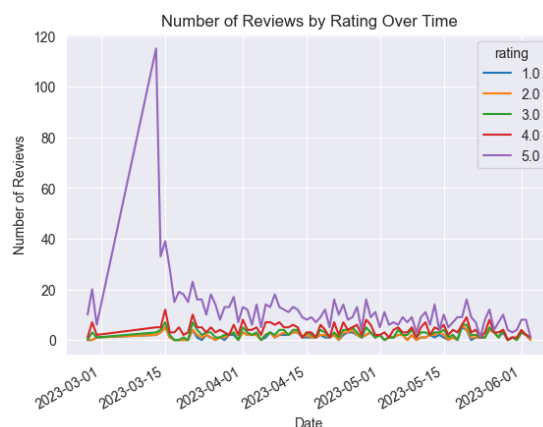


Fig. 10. Ilość opinii według oceny w czasie

Zgodnie z danymi dotyczącymi początku sprzedaży produktu, większość opinii użytkowników była pozytywna i wynosiła 5 gwiazdek. Można z tego wywnioskować, że popularność produktu rosła na samym początku jego wprowadzenia na rynek. Jednakże, wraz z upływem czasu, liczba pozytywnych opinii zaczęła spadać, a proporcja między pozytywnymi a negatywnymi opiniami ustabilizowała się na stałym poziomie.

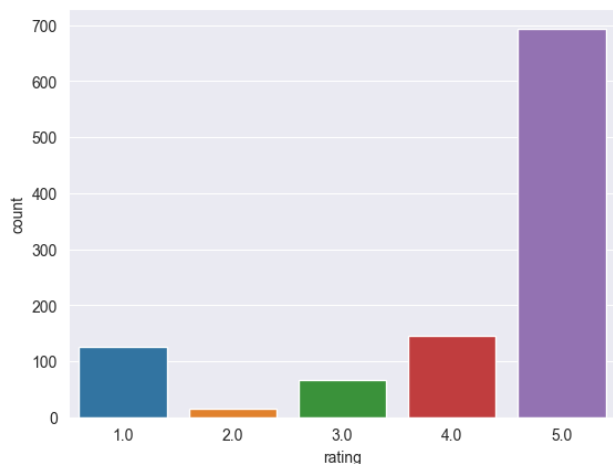


Fig. 11. Rozkład opinii

Można więc stwierdzić, że choć początkowo produkt cieszył się dużą popularnością i uznaniem, to jego sukces nie był długotrwały, a jego pozycja na rynku uległa stabilizacji. Warto jednak zauważyć, że mimo to produkt nadal zdobywa pozytywne opinie, co może świadczyć o jego wysokiej jakości oraz trwałości.

E. Treść

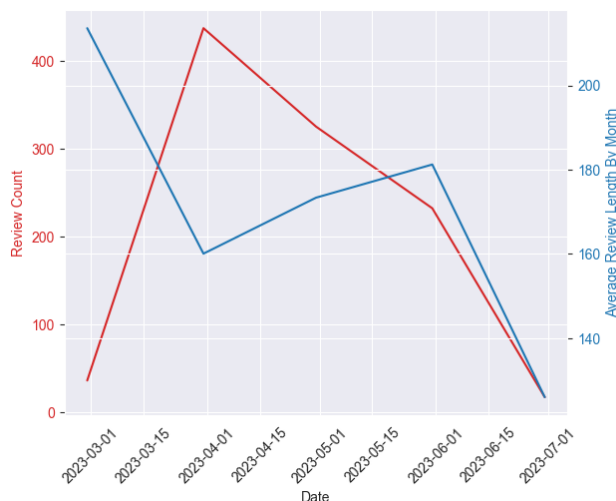


Fig. 12. Średnia długość recenzji według miesiąca i liczby recenzji

Analizując zależność pomiędzy liczbą recenzji a średnią długością, można stwierdzić, że wraz ze wzrostem liczby recenzji maleje średnia długość. Można to tłumaczyć tym, że kiedy produkt zyskuje na popularności i przyciąga więcej uwagi, pojawia się więcej osób, które chcą podzielić się swoją opinią na jego temat. Wśród tych osób mogą być zarówno osoby, które bardzo dobrze znają produkt i są w stanie przedstawić szczegółowe opinie na jego temat, jak i osoby, które dopiero co zapoznały się z nim i mają do powiedzenia tylko kilka słów.

F. Zależności pomiędzy danymi

Na podstawie tej korelacji możemy wyciągnąć kilka wniosków:

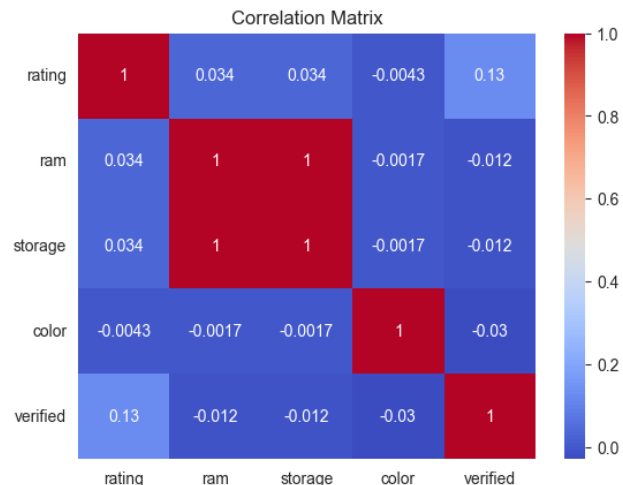


Fig. 13. Zależności pomiędzy danymi

1. Współczynnik korelacji między ratingiem a pozostałymi zmiennymi (ram, storage, color, verified) jest bardzo niski lub bliski zeru (co wynika z wartości współczynnika korelacji bliskich zero). Oznacza to, że brak jest silnej zależności między ratingiem a innymi zmiennymi, co sugeruje, że rating nie jest bezpośrednio związany z tymi cechami.

2. Współczynniki korelacji między ramem i storage'em oraz między tymi zmiennymi a pozostałymi (color, verified) są również bliskie zeru, co sugeruje brak zależności między nimi.

3. Współczynnik korelacji między weryfikacją a ratingiem jest stosunkowo wyższy niż dla pozostałych zmiennych, co może wskazywać na pewną zależność między nimi.

4. Współczynnik korelacji między kolorami a pozostałymi zmiennymi jest bliski zeru, co oznacza, że kolor nie jest silnie związany z innymi cechami.

Podsumowując, na podstawie tej korelacji można stwierdzić, że rating nie jest mocno skorelowany z innymi zmiennymi, takimi jak ram, storage czy kolor, ale istnieje pewna zależność między ratingiem a weryfikacją opinii.

4. EKSPERYMENTY

A. Analiza sentymentu

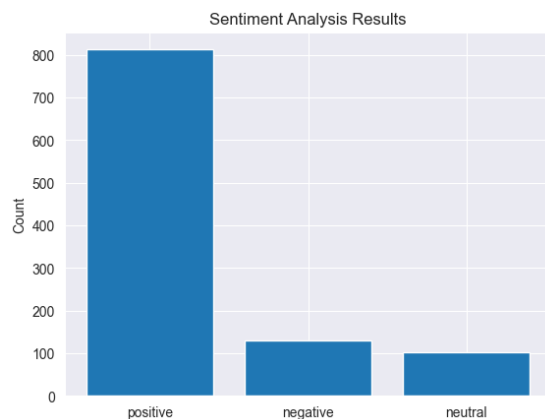


Fig. 14. Wynik analizy sentymentu modelem nltk

Pozytywny	Negatywny	Neutralny
813	131	103

Table 6. Wynik analizy sentymentu

B. Analiza negatywnych opinii

Po przeanalizowaniu wszystkich negatywnych opinii udało się zidentyfikować pięć głównych tematów:

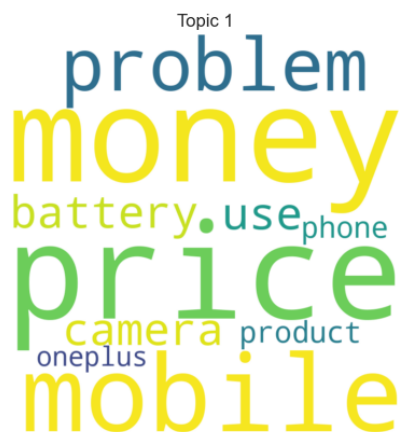


Fig. 15. Najpopularniejsze słowa tematu 1

Z tematu 1 wynika, że użytkownicy są zaniepokojeni ceną telefonów komórkowych, zużyciem baterii i wykorzystaniem aparatu. Oneplus jest popularną marką w tym temacie.



Fig. 16. Najpopularniejsze słowa tematu 2

Temat 2 sugeruje, że użytkownicy są niezadowoleni z jakości aparatu, jakości wyświetlacza i ogólnej jakości produktu swoich telefonów. Gry są również przedmiotem zainteresowania niektórych użytkowników.

W temacie 3 klienci omawiają doświadczenia związane z zakupami na Amazon, jakością usług i problemy, jakie napotkali ze swoimi telefonami Oneplus.

Temat 4 zwraca uwagę na kwestie związane z wydajnością i żywotnością baterii.

Wreszcie, temat 5 przynosi pozytywne komentarze na temat aparatu i żywotności baterii, ale wspomina również o negatywnych doświadczeniach z produktem i marką.



Fig. 17. Najpopularniejsze słowa tematu 3



Fig. 18. Najpopularniejsze słowa tematu 4



Fig. 19. Najpopularniejsze słowa tematu 5

Ogólnie, wydaje się, że klienci są przede wszystkim zainteresowani jakością aparatu, żywotnością baterii, jakością produktu i wydajnością. Obszary te mogą wymagać uwagi ze strony producentów w celu poprawy zadowolenia klientów.

Produkt z biegiem czasu wykazuje problemy z wydajnością oraz trwałością baterii, co potwierdza popularność 4 tematu na statystyce poniżej.

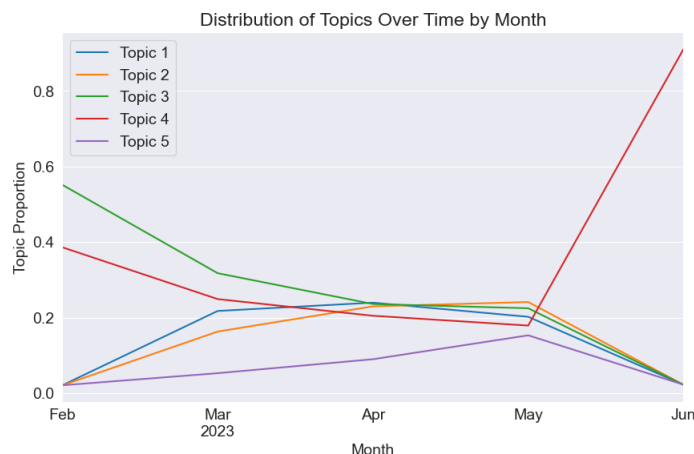


Fig. 20. Popularność tematów względem czasu

C. Analiza pozytywnych opinii

Po przeanalizowaniu wszystkich pozytywnych opinii udało się zidentyfikować pięć głównych tematów:

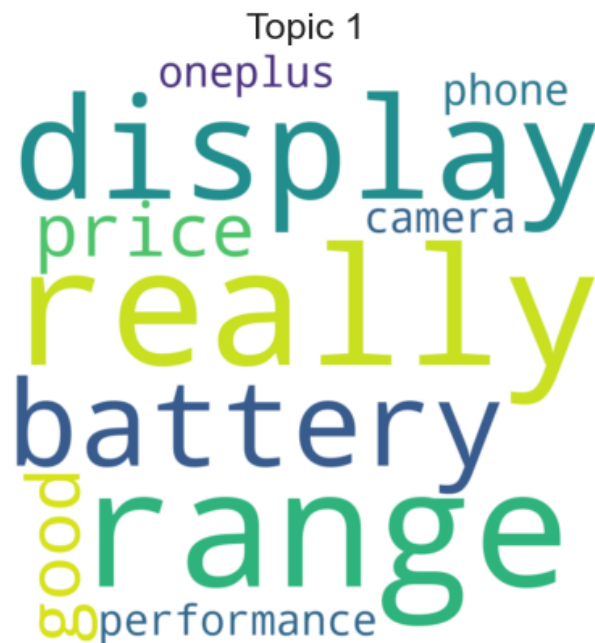


Fig. 21. Najpopularniejsze słowa tematu 1

Z tematu 1 możemy wywnioskować, że użytkownicy są zainteresowani zakresem funkcji i jakością ekranu telefonu OnePlus, a także jego żywotnością baterii i wydajnością. Użytkownicy wydają się również dyskutować na temat ceny telefonu i tego, jak odnosi się ona do jego ogólnej ceny.

W temacie 2 uwaga skupia się na ekranie telefonu i jego cenie, a użytkownicy szukają najlepszego urządzenia z dobrym czasem pracy baterii i funkcjami aparatu.

Temat 3 sugeruje, że użytkownicy cenią sobie telefon z doskonałym ekranem i długim czasem pracy baterii. Szukają również urządzenia mobilnego z doskonałymi funkcjami aparatu i są pod wrażeniem ogólnej jakości telefonu.

W temacie 4 użytkownicy wydają się być zaniepokojeni sto-

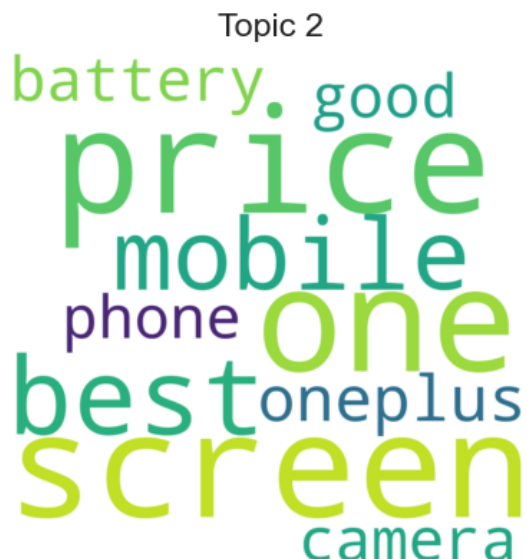


Fig. 22. Najpopularniejsze słowa tematu 2



Fig. 23. Najpopularniejsze słowa tematu 3

sunkiem jakości do ceny oferowanym przez markę OnePlus. Uważają, że produkt jest doskonałą inwestycją ze względu na ładny design i wydajność, które sugerują wysokiej jakości urządzenie mobilne.

Wreszcie, Temat 5 podkreśla znaczenie żywotności baterii, szybkiej wydajności i dobrych funkcji aparatu, aby użytkownicy uznali telefon za najlepszy.

Ogólnie, wydaje się, że użytkownicy są zainteresowani produktem ze względu na szereg funkcji i doskonałą wydajność. Istnieje jednak pole do poprawy w zakresie ceny produktu. Żywotność baterii, funkcje aparatu i jakość wyświetlania to najważniejsze priorytety dla użytkowników przy wyborze urządzenia mobilnego.



Fig. 24. Najpopularniejsze słowa tematu 4



Fig. 25. Najpopularniejsze słowa tematu 5

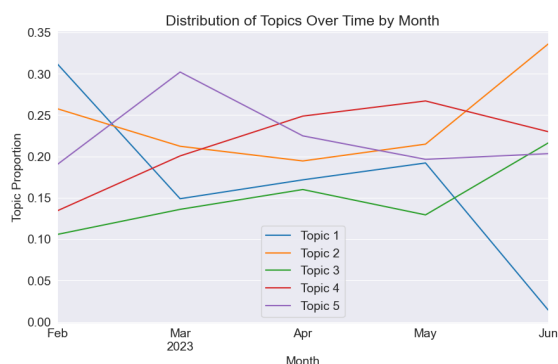


Fig. 26. Popularność tematów względem czasu

D. Analiza zależności sentymentu od oceny opinii

Statystyka chi-kwadrat wynosząca 9,81 i wartość p równa 0,007 wskazują, że istnieje statystycznie istotny związek między kategoriami nastrojów a zweryfikowanym statusem. Innymi słowy, rozkład kategorii sentymentu jest inny dla zweryfikowanych i niezweryfikowanych recenzji, a różnica ta prawdopodobnie nie wystąpiła przypadkowo.

Jedną z możliwych interpretacji tego wyniku jest to, że zweryfikowane recenzje częściej mają pozytywny sentyment niż recenzje niezweryfikowane. Należy jednak pamiętać, że korelacja nie oznacza związku przyczynowego i mogą istnieć inne czynniki, które wpływają zarówno na zweryfikowany status, jak i na sentyment.

Sentiment	False	True
Negative	25	106
Neutral	15	88
Positive	82	731

Table 7. Tabela krzyżowa wyników klasyfikacji sentymentu

Table 8. Wyniku testu chi-square

Chi-square statistic	p-value
9.813	0.0074

Statystyka chi-kwadrat wynosząca 9,813 i wartość p równa 0,0074 wskazują, że istnieje statystycznie istotny związek między kategoriami nastrojów a zweryfikowanym statusem. Rozkład kategorii sentymentu jest inny dla zweryfikowanych i niezweryfikowanych recenzji, a różnica ta prawdopodobnie nie wystąpiła przypadkowo.

Jedną z możliwych interpretacji tego wyniku jest to, że zweryfikowane recenzje częściej mają pozytywny sentyment niż recenzje niezweryfikowane. Należy jednak pamiętać, że korelacja nie oznacza związku przyczynowego i mogą istnieć inne czynniki, które wpływają zarówno na zweryfikowany status, jak i na sentyment.

E. Ekstrakcja cech i modele klasyfikacji

Została stworzona macierz cech TF-IDF. Następnie przeprowadzona eksploracja reguł asocjacyjnych.

Wyniki najistotniejsze względem siły zależności:

1. Istnieje silna zależność między posiadaniem telefonu koloru Sonic Black, oceną 3.0 a pojemnością 128 GB.
2. Klienci, którzy ocenili produkt na 2.0, częściej mieli urządzenie z pojemnością 128 GB i zweryfikowaną opinią.
3. Klienci preferujący telefony z pojemnością 256 GB i mają opinie z oceną 3.0 często wybierają kolor Galactic Silver.
4. Klienci, którzy posiadają telefony z pojemnością 256 GB, zweryfikowane informacje i ocenę 4.0, często wybierają model koloru Sonic Black.

5. Klienci, mające opinię z ratingiem 3.0 często mają telefon z pojemnością masową 128 GB.

Zostały przeprowadzone testy trzech modeli: regresji logistycznej, lasu losowego i gradient boosting do przewidywania ratingu. Każdy z modeli był trenowany na określonych danych dotyczących reguł i ratingu, które zostały podzielone na

zestawy treningowe i testowe. Celem testów było sprawdzenie skuteczności każdego z modeli w przewidywaniu wyników na zbiorze testowym. Do modeli regresji logicznej został użyty grid wyszukiwarka hiperparametrów. Do lasu losowego i gradientu boosting - losowa wyszukiwarka. Została wykorzystana metoda oversamplingu danych w celu przeciwdziałania nie zrównoważonym klasom w zbiorze treningowym.

Rating	Precision	Recall	F1-score	Support
1.0	0.68	0.65	0.67	26
2.0	0.00	0.00	0.00	4
3.0	0.13	0.15	0.14	13
4.0	0.30	0.23	0.26	35
5.0	0.80	0.86	0.83	131
Accuracy			0.67	209
Macro Avg	0.38	0.38	0.38	209
Weighted Avg	0.64	0.67	0.65	209

Table 9. Wyniki klasyfikacji modelu regresji logicznej

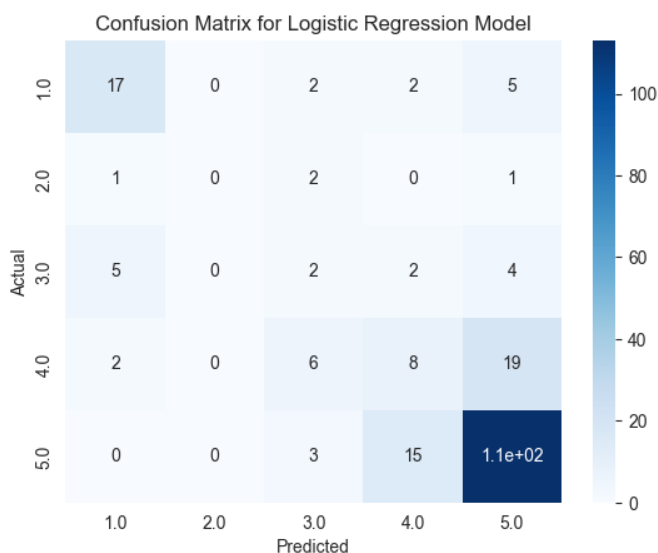


Fig. 27. Macierz pomyłek dla modelu regresji logicznej

1) Wyniki klasyfikacji Logistic Regression Model:

1. Precyzja (precision) dla większości klas jest umiarkowana, z wyjątkiem oceny 2.0, która ma precyzję równą 0.00. Precyzja wskazuje na zdolność modelu do poprawnego przewidywania pozytywnych przypadków.

2. Czułość (recall) dla większości klas jest również umiarkowana. Oceny 1.0 i 5.0 mają wyższą czułość, co oznacza, że model dobrze identyfikuje te oceny.

3. F1-score, które jest średnią ważoną precyzji i czułości, jest najwyższe dla oceny 5.0 i wynosi 0.83. Oznacza to, że model osiąga wysokie wyniki dla tej oceny.

4. Ogólna dokładność (accuracy) modelu wynosi 0.67, co oznacza, że około 67 procentów przypadków zostało poprawnie sklasyfikowanych.

Rating	Precision	Recall	F1-score	Support
1.0	0.69	0.69	0.69	26
2.0	0.00	0.00	0.00	4
3.0	0.17	0.08	0.11	13
4.0	0.44	0.23	0.30	35
5.0	0.78	0.95	0.86	131
Accuracy			0.72	209
Macro Avg	0.42	0.39	0.39	209
Weighted Avg	0.66	0.72	0.68	209

Table 10. Wyniki klasyfikacji modelu lasu losowego

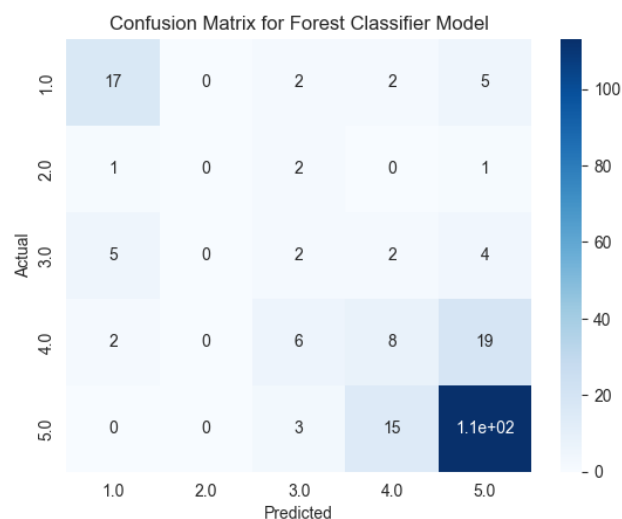


Fig. 28. Macierz pomyłek dla modelu lasu losowego

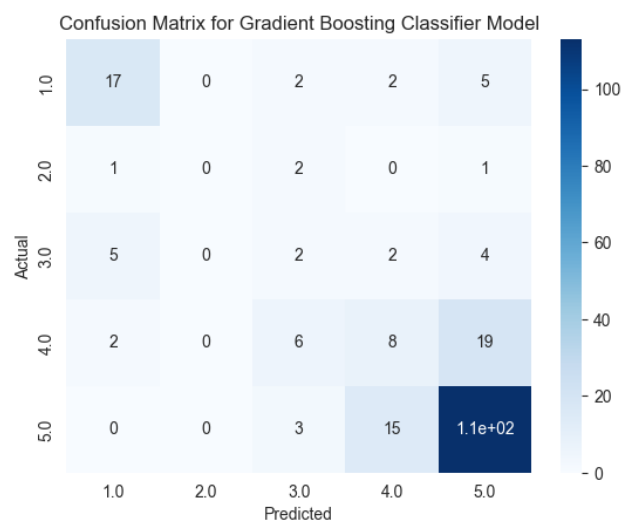


Fig. 29. Macierz pomyłek dla modelu gradient boosting

5. Na podstawie macierzy pomyłek (confusion matrix) można zauważyć, że model ma trudności z poprawnym rozpoznawaniem ocen 2.0 i 3.0.

Rating	Precision	Recall	F1-score	Support
1.0	0.70	0.54	0.61	26
2.0	0.00	0.00	0.00	4
3.0	0.12	0.15	0.13	13
4.0	0.33	0.34	0.34	35
5.0	0.80	0.83	0.82	131
Accuracy			0.66	209
Macro Avg	0.39	0.37	0.38	209
Weighted Avg	0.65	0.66	0.65	209

Table 11. Wyniki klasyfikacji modelu gradient boosting

2) Wyniki klasyfikacji RandomForestClassifier:

1. Precyzja dla większości ocen jest umiarkowanie wysoka, z wyjątkiem oceny 2.0, która ma precyzję równą 0.00.

2. Czułość dla większości ocen jest również umiarkowana. Ocena 5.0 ma najwyższą czułość, co wskazuje na dobrą zdolność modelu do identyfikowania tej oceny.

3. F1-score dla większości ocen jest stosunkowo wyższe w porównaniu do Logistic Regression Model. Najwyższe F1-score, wynoszące 0.86, osiągnięto dla oceny 5.0.

4. Ogólna dokładność modelu RandomForestClassifier wynosi 0.72, co jest nieco wyższe niż dokładność Logistic Regression Model.

5. Macierz pomyłek dla RandomForestClassifier jest identyczna jak w przypadku Logistic Regression Model, co sugeruje podobne wyniki klasyfikacji.

3) Wyniki klasyfikacji GradientBoostingClassifier:

1. Precyzja dla większości ocen jest umiarkowana. Ocena 2.0 ma precyzję równą 0.00. Czułość dla większości ocen jest również umiarkowana. Ocena 5.0 ma najwyższą czułość.

2. F1-score dla większości ocen jest zbliżone do wyników uzyskanych w Logistic Regression Model.

3. Ogólna dokładność modelu GradientBoostingClassifier wynosi 0.66, co jest nieco niższe niż w przypadku pozostałych dwóch modeli.

4. Macierz pomyłek dla GradientBoostingClassifier jest identyczna jak dla poprzednich modeli.

Podsumowując, RandomForestClassifier osiąga najwyższą ogólną dokładność (0.72) i najlepsze wyniki F1-score dla większości ocen, w szczególności dla oceny 5.0. Logistic Regression Model i GradientBoostingClassifier osiągają podobne wyniki, z nieco niższą ogólną dokładnością (0.67 dla Logistic Regression i 0.66 dla GradientBoostingClassifier). Wszystkie trzy modele mają trudności z poprawnym rozpoznawaniem oceny 2.0 i 3.0. To jest związane z tym, że jest niewystarczająca ilość danych odnośnie tych dwóch ocen.

F. Detekcja anomalii

Sprawdzone anomalii za pomocą algorytmu Isolation Forest.

Najpierw tekst jest przetwarzany za pomocą wektoryzacji TF-IDF, co pozwala na reprezentowanie tekstu jako liczb. Następnie wybierane są kolumny zawierające informacje o ratingu i sentymencie oraz konkatowane są ze sobą, aby utworzyć finalny zestaw danych.

Został poddany badaniu tydzień z największą ilością anomalii. Średni sentyment recenzji anomalii tego tygodnia jest

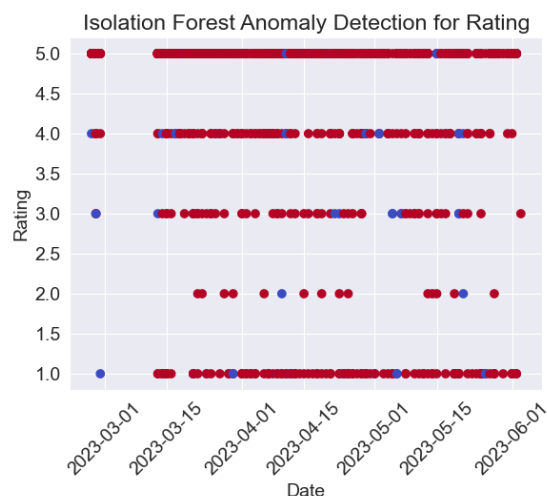


Fig. 30. Analiza anomalii względem oceny

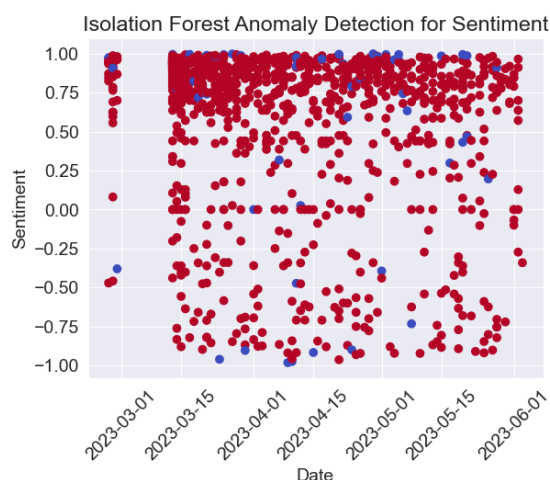


Fig. 31. Analiza anomalii względem sentymentu

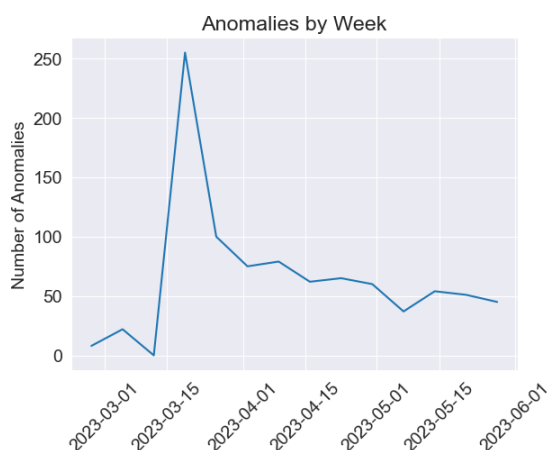


Fig. 32. Ilość anomalii względem tygodnia

pomiędzy neutralnym a pozytywnym.

Na podstawie dostarczonego wykrywania anomalii słów dla danego tematu możemy stwierdzić, że ludzie mają pozytywny

