

# Analiza zbioru danych "Patient-level dataset to study the effect of COVID-19 in people with Multiple Sclerosis" z wykorzystaniem nadzorowanego algorytmu uczenia maszynowego One-Class SVM(Support Vector Machines)

Anna Mrozek, Bartosz Panek i Przemysław Jura

## ARTICLE INFO

**Keywords:**  
One-Class SVM  
COVID-19

## STRESZCZENIE

Artykuł analizuje zbiór danych pacjentów ze stwardnieniem rozsianym (SM), którzy przeszli COVID-19, przy użyciu nadzorowanego algorytmu One-Class SVM. Celem jest identyfikacja przypadków o zwiększonym ryzyku ciężkiego przebiegu infekcji. Umożliwi to zidentyfikowanie charakterystycznych cech pacjentów bardziej narażonych na powikłania.

## 1. Wprowadzenie

Pandemia COVID-19 miała znaczący wpływ na osoby z chorobami przewlekłymi, w tym na pacjentów ze stwardnieniem rozsianym (SM). SM jest przewlekłą chorobą autoimmunologiczną, która wpływa na układ nerwowy i może prowadzić do trwałego uszkodzenia neuronów oraz ograniczenia funkcji motorycznych oraz poznawczych przez uszkodzenia mieliny oraz włókien nerwowych [1][2][3]. Osoby z MS są bardziej podatne na infekcje z powodu złożonego działania samej choroby, jej leczenia i naturalnego przebiegu [4]. Ze względu na charakter choroby i stosowane terapie immunosupresyjne, pacjenci z SM mogą być bardziej narażeni na cięższy przebieg infekcji wirusowych, w tym COVID-19.[5][6]

Analiza danych od pacjentów z SM, którzy przeszli COVID-19, może dostarczyć ważnych informacji na temat ryzyka powikłań i zidentyfikować czynniki przyczyniające się do poważniejszych objawów, co ostatecznie może pomóc w lepszym monitorowaniu i opiece nad tą grupą pacjentów.

## 2. Cel

Celem badania jest wykorzystanie nadzorowanego algorytmu One-Class SVM do analizy zbioru danych pacjentów ze stwardnieniem rozsianym, którzy przeszli COVID-19, w celu zidentyfikowania cech pacjentów o zwiększonym ryzyku ciężkiej infekcji. Oczekuje się, że wyniki analizy dostarczą wskazówek do opracowania bardziej ukierunkowanych strategii opieki i środków zapobiegawczych dla pacjentów ze stwardnieniem rozsianym w kontekście potencjalnych przyszłych pandemii lub innych zagrożeń wirusowych.

## 3. Przegląd literatury

W literaturze medycznej i naukowej, od początku pandemii COVID-19, wiele badań skupiało się na analizie wpływu wirusa SARS-CoV-2 na osoby cierpiące na choroby

przewlekłe i autoimmunologiczne, takie jak stwardnienie rozsiane (SM). Wyniki tych badań sugerują, że pacjenci z SM, zwłaszcza ci stosujący terapie immunosupresyjne, mogą być bardziej narażeni na cięższy przebieg COVID-19 oraz związane z nim powikłania[7][8]. U pacjentów z SM często obserwuje się obniżoną odporność oraz większą podatność na infekcje, co wynika zarówno z samej choroby, jak i z efektów leków tłumiących układ odpornościowy[9]. Dodatkowo w badaniach zwraca się uwagę na takie czynniki ryzyka jak wiek, płeć, rodzaj stosowanego leczenia oraz inne choroby towarzyszące, które mogą zwiększać ryzyko ciężkiego przebiegu infekcji u tej grupy pacjentów.

W ostatnich latach coraz większą popularność zyskuje zastosowanie algorytmów uczenia maszynowego, takich jak Support Vector Machines (SVM), do analizy danych medycznych. Algorytmy te pozwalają na wykrywanie wzorców i anomalii w dużych, zróżnicowanych zbiorach danych[10]. W szczególności algorytm One-Class SVM, używany do wykrywania anomalii, jest przydatny do identyfikacji przypadków wysokiego ryzyka w danych medycznych, gdzie przypadki nietypowe (np. pacjenci bardziej narażeni na ciężki przebieg COVID-19) występują rzadko. Badania pokazują, że One-Class SVM dobrze sprawdza się w analizie populacji pacjentów, gdy dostęp do dużych zbiorów danych o przypadkach zdrowych jest ograniczony, co często stanowi wyzwanie w analizie medycznej[11].

W kontekście badań nad COVID-19 i SM, algorytmy nadzorowanego uczenia maszynowego, takie jak SVM, umożliwiają dokładniejszą analizę czynników ryzyka oraz lepsze prognozowanie ciężkiego przebiegu choroby. Wyniki takich analiz mogą być przydatne nie tylko dla lekarzy, ale także dla decydentów zajmujących się zdrowiem publicznym, ponieważ pozwalają na opracowanie bardziej ukierunkowanych strategii opieki dla osób z SM, które są bardziej narażone na zagrożenia związane z wirusami, takimi jak COVID-19.

## 4. Metodologia

W badaniu zastosowano algorytm nadzorowanego uczenia maszynowego One-Class SVM w celu identyfikacji pacjentów z podwyższonym ryzykiem ciężkiego przebiegu COVID-19 w grupie osób ze stwardnieniem rozsianym. Analiza obejmowała cechy kliniczne pacjentów, takie jak wiek, płeć, typ leczenia oraz choroby współistniejące, aby zidentyfikować wzorce związane z większą podatnością na powikłania.

### 4.1. Dataset

Z inicjatywy MS Global Data Sharing Initiative (GDSI) zbadano, jak leki immunosupresyjne lub immunomodulujące wpływają na COVID-19 i jego przebieg u osób z MS. GDSI miała na celu zwiększenie skali zbierania danych dotyczących COVID-19 i dostarczenie społeczności związanej z MS informacji opartych na danych podczas pandemii [12]. W ramach GDSI wybrano kluczowe zmienne obejmujące informacje o COVID-19, stopniu jego ciężkości, leczeniu, dane demograficzne, historię i nasilenie MS, stosowanie leków modyfikujących przebieg choroby (DMT), choroby współistniejące i wybrane zachowania związane ze stylem życia, takie jak palenie tytoniu. Globalna społeczność MS współpracowała, przekazując dokumentację statusu COVID-19 u osób z MS za pośrednictwem centralnej platformy udostępnionej przez QMENTA [13].

Ten zbiór danych został zebrany za pomocą narzędzia do szybkiego wprowadzania danych, które umożliwiło klinicystom, osobom ze stwardnieniem rozsianym (PwMS) lub ich przedstawicielom bezpośrednie wprowadzanie informacji do centralnej platformy GDSI. Narzędzie to zawierało kwestionariusz oparty na wcześniej ustalonych zmiennych i nie gromadziło bezpośrednich danych osobowych, aby chronić prywatność użytkowników. Narzędzie zostało wyłączone 3 lutego 2022 roku.

Zbiór danych obejmuje informacje o 1141 osobach ze stwardnieniem rozsianym (PwMS). Aby zapewnić zgodność danych z wytycznymi HIPAA, przeprowadzono proces deidentyfikacji.

Ponieważ w danych nie zbierano imion pacjentów, deidentyfikacja skupiła się na datach i wieku pacjentów. Daty w kolumnie „stop-or-end-date-combined” zostały przesunięte o losową liczbę dni (między -15 a 15), aby uniemożliwić identyfikację na podstawie dat. Wiek pacjentów sklasyfikowano w cztery grupy: 0 dla osób w wieku 0–17 lat, 1 dla osób między 18 a 50 lat, 2 dla osób między 51 a 70 lat, oraz 3 dla osób powyżej 71 lat. Dzięki temu żadne dokładne wartości wieku powyżej 90 lat nie zostały ujawnione. Po klasyfikacji zmiennych i wdrożeniu odpowiednich środków ostrożności dane zostały zdeidentyfikowane i spełniają standardy HIPAA, zachowując jednocześnie wartość badawczą. Ponadto, aby zapewnić ochronę prywatności, zastosowano techniki takie jak Kanonimizacja oraz różnorodność.

Zbiór danych obejmuje zestaw z góry określonych zmiennych, takich jak płeć, kategoria wiekowa, typ MS, wynik EDSS, status palenia oraz kategoria BMI. Te zmienne

dostarczają informacji o demografii pacjentów, ich stanie klinicznym oraz symptomach związanych z COVID-19.

### 4.2. Opis metody

Jednoklasowy SVM (One-Class Support Vector Machine, OCSVM) to algorytm przeznaczony do wykrywania anomalii w zbiorze danych. Zamiast klasyfikować dane do dwóch lub więcej klas, jak w klasycznym SVM, OCSVM koncentruje się wyłącznie na danych normalnych.[14][15] Jego celem jest zbudowanie granicy wokół większości przypadków normalnych, tworząc "strefę normalności", która odróżnia standardowe przypadki od anomalii.

Główna zasada działania OCSVM opiera się na maksymalizacji marginesu między przypadkami normalnymi a granicą, która oddziela normę od anomalii. Dzięki temu model lepiej rozpoznaje dane odstające, które znajdują się poza wyznaczoną strefą.[16][17]

Kluczowym elementem działania OCSVM jest funkcja decyzyjna oparta na hiperpłaszczyźnie, która oddziela dane od początku układu współrzędnych. Wzór na hiperpłaszczyznę wyrażany jest jako:

$$f(x) = \mathbf{w} \cdot \mathbf{x} - \rho \quad (1)$$

gdzie:

- $\mathbf{w}$  – wektor normalny do hiperpłaszczyzny, wyznaczony podczas procesu uczenia,
- $\mathbf{x}$  – próbka danych,
- $\rho$  – jest wartością progową (bias), która definiuje margines.

OCSVM posiada szereg parametrów, dzięki którym możliwe jest dostosowywanie modelu do specyficznych potrzeb:

- **Parametr  $\nu$ :** Kontroluje liczbę obserwacji uznawanych za anomalie. Parametr ten przyjmuje wartości w przedziale od 0 do 1 i determinuje udział anomalii, jakie model może tolerować w zbiorze treningowym (np. jeśli  $\nu = 0.05$ , oznacza to, że około 5% próbek zostanie sklasyfikowanych jako anomalie).
- **Parametr  $\gamma$ :** Wpływa na kształt granicy decyzyjnej poprzez dopasowywanie modelu do danych. Mniejsza wartość  $\gamma$  oznacza „szerszą” granicę, obejmującą więcej punktów, natomiast większa wartość  $\gamma$  powoduje, że granica jest bardziej dopasowana do danych, co może zwiększać ryzyko nadmiernego dopasowania (overfitting).
- **Jądro (*ang. kernel*):** One-Class SVM często wykorzystuje funkcje jądrowe, aby modelować nieliniowe granice decyzyjne i lepiej dopasować się do danych, które nie są liniowo rozdzielne w przestrzeni cech. Wybór jądra, zwany również „sztuczką jądra”, wpływa na charakter dopasowania modelu i jego wydajność.

Podczas treningu OCSVM analizuje wy cznie dane normalne, co sprawia,  e jest szczególnie u yteczny w sytuacjach, gdy anomalie s  rzadkie lub trudne do zidentyfikowania. Model staje si  wtedy bardziej niezawodny w rzeczywistych zastosowaniach, takich jak wykrywanie oszustw, monitorowanie awarii czy zabezpieczanie sieci komputerowych. OCSVM mo e wykorzystac r  ne funkcje j dra, co umo liwia mu wykrywanie zar wno prostych, jak i bardziej z ozonych odchyle . [18]

### 4.3. Opis przeprowadzonych oblicze 

1. Przygotowanie i czyszczenie danych - rozpocz liśmy od wczytania danych oraz wybrania interesuj cych nas cech, kt re nast pnie zosta y przekszta cone w spos b umo liwiaj cy przeprowadzenie oblicze :

- Dla kolumn binarnych (yes/no) warto ci zosta y zamienione na liczby 0 i 1.
- Dla kolumn kategorycznych i ordinalnych (np. age-in-cat, covid19-outcome-levels-2, report-source) zosta y przypisane warto ci liczbowe, zgodnie z przyj tym mapowaniem.
- Brakuj ce warto ci w danych zosta y wype nione zerami, co pozwoli o na unikni cie problem w podczas analizy i treningu modelu.

2. Tworzenie nowych zmiennych: symptom-score oraz comorbidity-score, aby skonsolidowa  informacje o objawach i chorobach wsp listniej cych.

- Stworzy lmy dwie kolumny symptom-score oraz comorbidity-score, kt re zliczaj  ilo   objaw w wirusa COVID-19 oraz liczby chor b wsp listniej cych dla ka dego pacjenta.
- Kolumny te by y nast pnie skalowane przy u yciu StandardScaler, co pozwala na lepsz  interpretacj  wyników oraz standaryzacj  w zakresie modelowania.

Efekt m by o, uzyskanie warto ci numerycznych, kt re reprezentuj  intensywno   symptom w oraz liczb  chor b wsp listniej cych dla ka dego pacjenta.

3. Dane zosta y podzielone na zbi r treningowy ( $X_{train}$ ) i testowy ( $X_{test}$ ) za pomoc  funkcji `train_test_split` z biblioteki `sklearn`, przy czym 70% danych przydzielono do zbioru treningowego. Zbi r treningowy s u y do uczenia modelu, natomiast testowy zosta  rozszerzony do wszystkich danych, wi c zawiera  w sobie nowe nieznanne dane dla modelu, jak i poprzednie dane. Wynika o to z l twiejszego testowania danych, a wyniki i wnioski pozostaw y podobne, a nawet lepsze po obserwacji, wynika o to mi dzy innymi z wi kszej pr by. Ustawienie argumentu `random_state=42` zapewnia reprodukowalno   wyników, dzi ki czemu podzia  danych jest zawsze taki sam przy kolejnych uruchomieniach. Warto zaznaczy , i  na potrzeby dalszej oceny modelu pacj ci hospitalizowani zostali uznani za anomalie poniewa  najci iej przeszli chorob  (ok. 2%). Dlatego do nauki modelu zostali wykorzystani jedynie pacj ci nie hospitalizowani, w celu nauczania danych na danych normalnych.

4. Trening modelu One-Class SVM - zastosowa lmy algorytm One-Class Support Vector Machine (SVM) z j drem `rbf`, kt ry dobrze radzi sobie z modelowaniem nieliniowych granic decyzyjnych. Model trenowano, aby rozpoznawa  typowe cechy pacjent w COVID-19 i SM, a nast pnie by przewidywa , kt re pr bki s  podobne do tego wzorca, a kt re mog  by  anomaliami. Po treningu model przypisywa  ka dej pr bce etykiet  1 (normalna) lub -1 (anomalia). Pr bki oznaczone jako anomalie zawarto w nowej ramce `anomalies` w celu dalszej analizy. Model trenowano na r  nych parametrach, czy to wielko ci przyj tych anomalii, kszta tu granicy decyzyjnej czy w sp czynnika j dra. Najlepsze wyniki osi ga  dla przyj tych 2% anomalii i w sp czynnika j dra `'scale'`.

5. Analiza wyników - na podstawie wykrytych anomalii przeprowadzili my szereg analiz, tworz c wizualizacje, kt re pomagaj  zrozumie  charakterystyk  anomalii [19]

### 4.4. Wykorzystane metryki oceny

**Raport klasyfikacji (ang. classification report)** to podsumowanie wyników modelu klasyfikacyjnego, kt re przedstawia kluczowe metryki takie jak precyzja (`precision`), czu o   (`recall`), i F1-score dla ka dej klasy. Precyzja wskazuje, jak wiele z klasyfikacji pozytywnych jest poprawnych, podczas gdy czu o   pokazuje, ile przyk ad w rzeczywi cie pozytywnych zosta o poprawnie sklasyfikowanych. F1-score jest s redni  harmoniczn  precyzji i czu o ci, co sprawia,  e dobrze nadaje si  do oceny modeli na danych niezbalansowanych. Raport zawiera r wnie  wska nik wsparcia (`support`), kt ry pokazuje liczb  przyk ad w w ka dej klasie.

**Dok adno  **, inaczej zwana `accuracy`, to stosunek liczby poprawnych predykcji (zar wno pozytywnych, jak i negatywnych) do ca kowitej liczby przyk ad w w danych testowych. Jest to og lny wska nik skuteczno ci modelu. Jednak w przypadku danych niezbalansowanych, dok adno   mo e by  myl ca, poniewa  wysoka warto   mo e wynika  z dominacji jednej klasy. Dlatego warto uzupe nia  analiz  dok adno ci bardziej szczeg lowymi metrykami, jak np. F1-score czy `False Positive Rate`.

**Funkcja decyzyjna (ang. decision\_function)** w modelach takich jak One-Class SVM okre la odleg o   przyk ad w od granicy decyzyjnej modelu. Wyniki funkcji decyzyjnej wskazuj , na ile pewnie model klasyfikuje przyk ad jako normalny lub anomalie. Warto ci dodatnie sugeruj ,  e przyk ad jest wewn trz granicy klasy normalnej, a warto ci ujemne sugeruj ,  e przyk ad jest anomalie. Warto ci bli zej zera oznaczaj ,  e przyk ad znajduje si  blisko granicy decyzyjnej.

**Odsetek anomalii** okre la, jaka cz   danych zosta a sklasyfikowana jako anomalie przez model. W modelach takich jak One-Class SVM parametr `nu` kontroluje, jaki procent danych jest uwa any za anomalie podczas trenowania modelu. Je li odsetek anomalii w wynikach jest znacznie wi kszy ni  zak adany, mo e to oznacza  nadmiern  surowo   modelu. Z kolei zbyt niski odsetek mo e sugerowa ,  e model nie wykrywa rzeczywistych anomalii.

**False Positive Rate (FPR)** to wskaźnik fałszywych alarmów, obliczany jako stosunek liczby fałszywie pozytywnych klasyfikacji (FP) do liczby wszystkich rzeczywistych przykładów negatywnych (FP + TN). Wysoki FPR oznacza, że model często błędnie klasyfikuje normalne przypadki jako anomalie, co może być problematyczne w zastosowaniach wymagających dużej precyzji (np. w diagnostyce medycznej). Zmniejszenie FPR można osiągnąć przez odpowiednie dostrojenie parametrów modelu lub zmianę progu decyzyjnego.

## 5. Wyniki

### 5.1. Ocena modelu One-Class SVM

	precision	recall	f1-score	support
0	0.22	0.53	0.31	15
1	0.99	0.97	0.98	1126
accuracy			0.97	1141
macro avg	0.60	0.75	0.65	1141
weighted avg	0.98	0.97	0.97	1141

Dokładność: 96.84%  
Funkcja decyzyjna (Średni score): 0.34  
Dokładność dla założonych anomalii: 8/15  
Odsetek anomalii: 3.24%  
False Positive Rate (FPR): 0.47

Figure 1: Ocena modelu One-class SVM

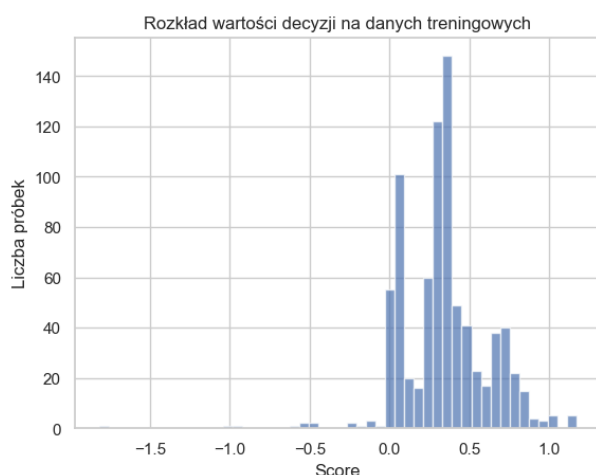


Figure 2: Rozkład wartości decyzji dla danych modelu One-class SVM

Na potrzeby oceny modelu, zostały wyodrębnione potencjalne dane które prawdopodobnie były anomalią to znaczy przypadkami, które wymagały hospitalizacji i poważniejszego leczenia. To właśnie ta grupa osób została uznana za tą która najczęściej przeszła chorobę. Za potencjalne anomalie zostały uznane przypadki, które w danych w kolumnie covid19\_admission\_hospital zawierały informacje o hospitalizacji. Dodatkowo fakt, iż

**jedynie ok 2% osób badanych trafiło do szpitala potwierdziło założenie, że to właśnie tam należy doszukiwać się potencjalnych anomalii.**

Model klasyfikacyjny został oceniony na podstawie kilku metryk, takich jak precyzja (precision), czułość (recall), F1-score, oraz wsparcie (support). Dla klasy 0 (anomia) precyzja wyniosła 0.22, co oznacza, że jedynie 22% przypadków sklasyfikowanych jako anomalie było poprawnych. Czułość wyniosła 0.53, co wskazuje, że model wykrył 53% rzeczywistych anomalii. Niski F1-score (0.31) podkreśla trudności modelu w skutecznym identyfikowaniu anomalii, co jest szczególnie problematyczne przy tylko 15 przykładach rzeczywistych anomalii (support). Dla klasy 1 (normalny przypadek), wyniki są bardzo wysokie – precyzja 0.99 i czułość 0.97 świadczą o dużej skuteczności modelu w identyfikacji normalnych przypadków, co znajduje odzwierciedlenie w wysokim F1-score równym 0.98. Średnia dokładność modelu (accuracy) wyniosła 96.84%, co sugeruje ogólną skuteczność klasyfikacji, ale maskuje problem niskiej wydajności w detekcji anomalii.

Odsetek anomalii w zbiorze danych wynosi 3.24%, co wskazuje na dużą nierównowagę klas. False Positive Rate (FPR) na poziomie 0.47 oznacza, że model często błędnie klasyfikuje normalne przypadki jako anomalie. Dodatkowo, średni wynik na danych testowych z funkcji decyzyjnej wyniósł 0.34, co wskazuje, że większość przykładów testowych znajduje się blisko granicy decyzyjnej modelu. To może sugerować, że model ma trudności z wyraźnym rozdzieleniem klas, szczególnie w przypadku anomalii.

### 5.2. Wynik analizy

Wykresy znajdują się na końcu opracowania.

## 6. Opis wyników

Opis wykresów znajdujących się na końcu opracowania:

Figure 3: Na wykresie widzimy histogram z liczbą próbek oznaczonych jako normalne i anomalia przez model One-Class SVM. Z wykresu wynika, że większość próbek została sklasyfikowana jako normalne (oznaczone jako 1), podczas gdy mniejsza liczba próbek została uznana za anomalie. Wysoka liczba normalnych próbek w stosunku do anomalii wskazuje, że model wykrywa anomalie rzadziej, co jest zgodne z celem detekcji anomalii, ponieważ anomalie rzadziej występują.

Figure 4: Histogram pokazuje, że większość przypadków anomalii ma niski symptom\_score, głównie na poziomie 0, z kilkoma przypadkami rozproszonymi na wyższych wartościach aż do 10. Niskie wartości mogą sugerować, że większość anomalii wykazuje niewiele lub żadnych symptomów, choć istnieją wyjątki, gdzie symptom score jest wyższy.

Figure 5: Histogram przedstawia, że większość anomalii ma comorbidity\_score równy zero, co oznacza brak chorób współistniejących, chociaż kilka przypadków posiada wyższe wartości aż do 4. To sugeruje, że choć wiele

przypadkōw anomalii nie ma dodatkowych schorzeń, niektōre z nich charakteryzujā siē złoŹonymi chorobami wspōlistniejācymi.

Figure 6: Histogram pokazuje rozkłād Symptom\_score dla normalnych przypadkōw. Najwiēcej przypadkōw ma symptom\_score wynoszący 0, ale przypadki sā bardziej równomiernie rozproszone w przedziale od 1 do 7. Sugeruje to, Źe normalne przypadki mogā mieć rōżny poziom symptomōw, ale głōwnie wyrōżniają siē skrajnie niskimi wartościami.

Figure 7: Tutaj widzimy, Źe wiēkszość przypadkōw normalnych ma comorbidity\_score bliski zero, choć niektōre przypadki osiagajā wartośc 1. Wskazuje to, Źe wiēkszość przypadkōw normalnych nie ma chorōb wspōlistniejācych, ale mogā występować łagodne wspōlistniejāce schorzenia.

Figure 8: Na wykresie przedstawiono średnie wartościcech binarnych dla przypadkōw zaklasyfikowanych jako anomalie. Sā to cechy, ktōre wskazujā obecnośc lub brak pewnych objawōw lub stanōw zdrowotnych. Wysokie średnie wartościce wskazujā na czēstsze występowanie danej cechy wśrōd anomalii: Najwiēksze średnie wartościce dotyczą cech takich jak current\_dmt (obecnie przyjmowana terapia modyfikujāca chorobę), covid19\_self\_isolation (izolacja w zwiāzku z COVID-19), oraz covid19\_has\_symptoms (obecnośc symptomōw COVID-19). Oznacza to, Źe osoby klasyfikowane jako anomalie czēsto sā poddane leczeniu, majā objawy lub byli w izolacji. Kolejne wysokie cechy to np. has\_comorbidities (choroby wspōlistniejāce) oraz rōżne objawy zwiāzane z COVID-19, takie jak dry\_cough (suchy kaszel) i fatigue (zmęczenie), co sugeruje, Źe anomalie mogā być powiāzane z przypadkami COVID-19 lub innymi powaŹnymi objawami chorobowymi.

Figure 9: Na tym wykresie poziomym przedstawiono średnie wartościce dla rōżnych cech binarnych w normalnych przypadkach (bez anomalii). Cecha o najwiēkszej średniej wartościce to current\_dmt, co wskazuje, Źe pacjenci normalni czēsto uŹywajā terapii modyfikujācych chorobę (DMT). Inne czēsto występujāce cechy to current\_or\_former\_smoker oraz covid19\_self\_isolation, oraz covid19\_has\_symptoms, co sugeruje, Źe objawy COVID-19 i byli bādź obecni palacze t znaczna grupa osōb w normalnej grupie.

Figure 10: Wiēkszość anomalii dotyczy osōb młodszych (kategoria 1), co moŹe sugerować, Źe młodsze osoby majā wiēksze ryzyko występowania anomalii w danych.

Figure 11: Zdecydowana wiēkszość przypadkōw anomalii dotyczy osōb z niŹszym BMI (kategoria 0), co moŹe wskazywać na zwiāzek miēdzy BMI a nietypowymi przypadkami.

Figure 12: Wśrōd anomalii przeważajā kobiety (kategoria 1), co sugeruje, Źe w danych kobiet czēściej występujā anomalie.

Figure 13: Wiēkszość anomalii dotyczy osōb bez przypisania do ms\_type2 (kategoria 0), ale niektōre przypadki naleŹą do kategorii 1 lub 2.

Figure 14: Rozkłād wskazuje, Źe wśrōd anomalii dominujā łagodniejsze wyniki COVID-19 (kategoria 0), choć sā teŹ przypadki umiarkowane (kategoria 1) i cięŹkie (kategoria 2).

## 7. Podsumowanie

Zgromadzony zbiōr danych liczył jedynie ponad 1100 przypadkōw, co w kontekście analiz statystycznych i uczenia maszynowego jest liczbā stosunkowo niewielkā, szczegōlnie przy analizie złoŹonych zjawisk, takich jak przebieg COVID-19 w określonych populacjach. Dodatkowo dane te charakteryzowały siē duŹā iloścā brakujācych wartościce – aŹ okołō 70% przypadkōw zawierało znaczne luki w informacjach, co ograniczało ich uŹytecznościce. Tylko 2% przypadkōw dotyczyło osōb hospitalizowanych, co dodatkowo zmniejszało wartościce predykcijnā modelu, poniewā cięŹkie przypadki stanowiły bardzo niewielkā czēść danych i byli słabo reprezentowane.

Problemem była rōwnieŹ nierównomiernościce rozkłādu danych w kluczowych grupach, takich jak płeć i wiek, co prowadziło do stronniczościce wyników. Na przykłād, nierównomierna liczba kobiet i mēŹczyzn w próbce lub rōŹnice w liczbie obserwacji pomiēdzy grupami wiekowymi powodowały, Źe model mōgł błędnie wyciagāc wnioski na podstawie nadreprezentowanych cech. Brak równowagi w danych utrudniał takŹe identyfikacjē rzeczywistych wzorcōw w populacji i prowadził do fałszywego postrzegania wiarygodnościce wyników.

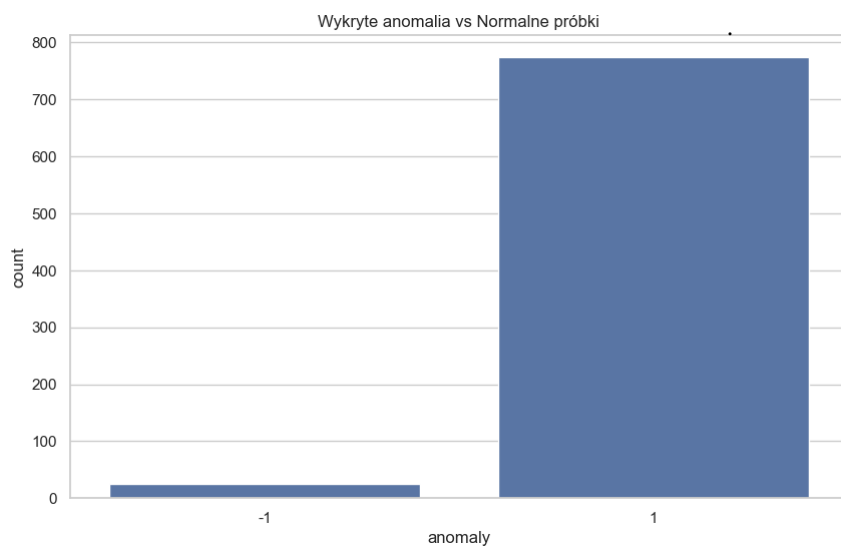
Dodatkowym wyzwaniem w analizie takich danych jest specyfika COVID-19 jako choroby o niejasnych i zrōŹnicowanych objawach klinicznych. COVID-19 moŹe manifestować siē w bardzo szerokim spektrum od całkowicie bezobjawowego przebiegu, poprzez łagodne symptomy, aŹ po cięŹkie przypadki wymagajāce hospitalizacji i intensywnej terapii. Ta zmiennościce znacząco utrudnia stworzenie jednorodnych wzorcōw klasyfikacyjnych, szczegōlnie gdy dane sā ograniczone i niekompletne. W przypadku omawianego zbioru brak dostatecznej reprezentacji cięŹkich przypadkōw oraz brak dokłādnych informacji o objawach dodatkowo ograniczał moŹliwościce identyfikacji kluczowych cech rōŹnicujācych normalne przypadki od anomalii.

W tych okolicznościcach zastosowanie modelu one-class SVM nie było włāściwym podejściem. Model ten zakłāda, Źe dostępane dane normalne sā wystarczajāco reprezentatywne, by umoŹliwić wykrycie anomalii. JednakŹe w tym przypadku, z powodu małej liczby przypadkōw, braku wystarczajācych przykłādōw rzeczywistych anomalii (np. cięŹkich przypadkōw COVID-19), a takŹe brakōw danych i ich nierównomiernego rozkłādu, model nie był w stanie skutecznie odrōŹnić anomalii od danych normalnych. W konsekwencji one-class SVM wykazał fałszywie wyniki.

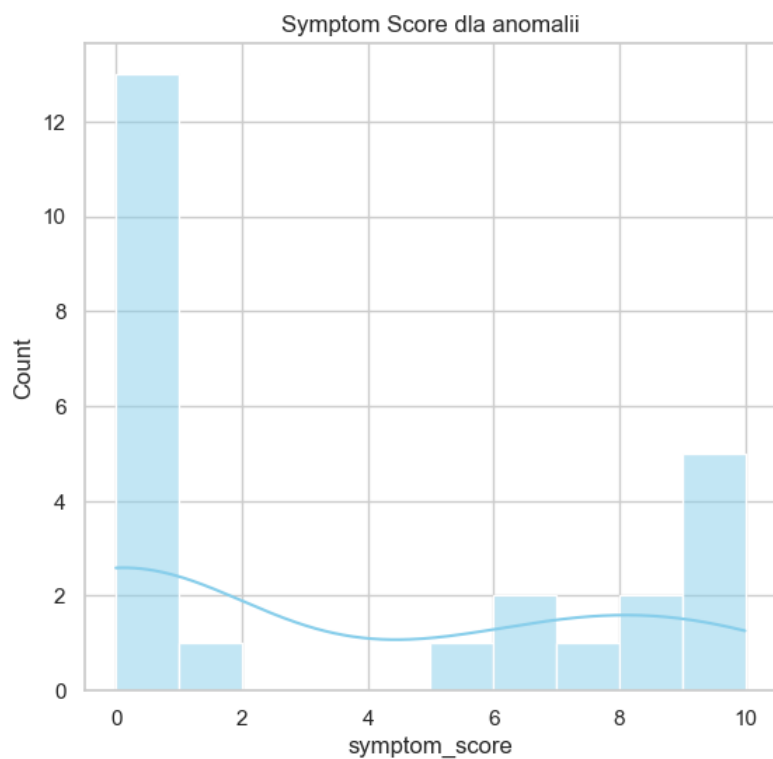
**Podsumowujāc, niedostateczna liczba przypadkōw, liczne braki danych, nierównomiernościce ich rozkłādu oraz brak reprezentatywnościce cięŹkich przypadkōw sprawiły, Źe analiza oparta na one-class SVM nie była miarodajna. Wyniki uzyskane przez model byli w duŹej mierze fałszywie pozytywne i nie miały praktycznego znaczenia. Aby przeprowadzić skuteczniejszā analizę, konieczne byloby zgromadzenie wiēkszej, bardziej zrōwnowaŹonej i kompletnej próbki danych, uwzględniajācej wiēkszā liczbę przypadkōw klinicznych i cięŹkich przypadkōw COVID-19.**

## 8. Bibliografia

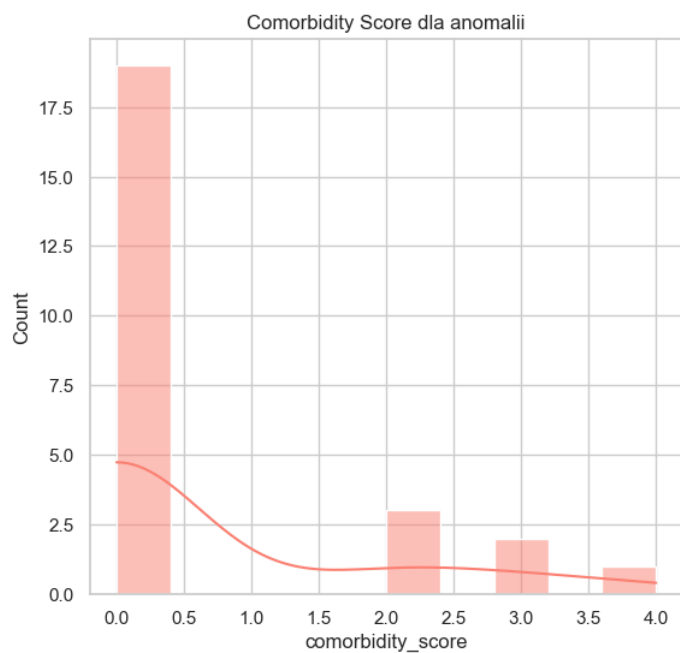
- [1] Diag.pl. (n.d.). Jakie s pierwsze objawy i sposoby leczenia stwardnienia rozsianego. Diag.pl. <https://diag.pl/pacjent/artykuly/jakie-sa-pierwsze-ob-jawy-i-sposoby-leczenia-stwardnienia-rozsianego/>
- [2] Medicover. (n.d.). Stwardnienie rozsiane – objawy, przyczyny i leczenie. Medicover. <https://www.medicover.pl/ozdrowiu/stwardnienie-rozsiane-objawy-przyczyny-i-leczenie-6619,n,192>
- [3] Calabresi PA. Diagnosis and management of multiple sclerosis. *Am Fam Physician*. 2004 Nov
- [4] Montgomery S, Hillert J, Bahmanyar S. Hospital admission due to infections in multiple sclerosis patients. *Eur J Neurol*. 2013 Aug
- [5] Wikipedia. (2025). Stwardnienie rozsiane. Wikipedia. [https://pl.wikipedia.org/wiki/Stwardnienie\\_rozsiane](https://pl.wikipedia.org/wiki/Stwardnienie_rozsiane)
- [6] Polskie Towarzystwo Stwardnienia Rozsianego. (2020). COVID-19 a SM. PTSR. <https://ptsr.org.pl/strona/133,covid-19-a-sm>
- [7] Sormani MP, De Rossi N, Schiavetti I, Carmisciano L, Cordioli C, Radaelli M, et al. Disease modifying therapies and COVID-19 severity in multiple sclerosis. *Ann Neurol*. 2021 Apr
- [8] Louapre C, Collongues N, Stankoff B, Giannesini C, Papeix C, Bensa C, et al. Clinical characteristics and outcomes in patients with coronavirus disease 2019 and multiple sclerosis. *JAMA Neurol*. 2020 Sep
- [9] Simpson-Yap S, De Brouwer E, Kalincik T, Rijke N, Hillert JA, Walton C, et al. Associations of Disease-Modifying Therapies With COVID-19 Severity in Multiple Sclerosis. *Neurology*. 2021 Nov
- [10] Schiff MA, Rae-Grant A, Gilden D, Franklin GM. Practice guideline: Disease-modifying therapies for adults with multiple sclerosis: Report of the guideline development, dissemination, and implementation subcommittee of the American Academy of Neurology. *Neurology*. 2019 Jan
- [11] Erfani P, Mitchell AJ, Hameed S, Heydarpour P, Ghaffaripour R, Sahraian MA. Systematic review of health-related quality of life in multiple sclerosis patients: The impact of pharmacological treatments and lifestyle. *J Neurol Sci*. 2016 Dec
- [12] Peeters LM, Parciak T, Walton C, Geys L, Moreau Y, De Brouwer E, et al. COVID-19 in people with multiple sclerosis: A global data sharing initiative. *Mult Scler Houndmills Basingstoke Engl*. 2020 Sep
- [13] Simpson-Yap S, De Brouwer E, Kalincik T, Rijke N, Hillert JA, Walton C, et al. Associations of Disease-Modifying Therapies With COVID-19 Severity in Multiple Sclerosis. *Neurology*. 2021 Nov
- [14] Physionet. (2025). Patient-level data for COVID and MS. Physionet. <https://physionet.org/content/patient-level-data-co-vid-ms/1.0.1/>
- [15] GeeksforGeeks. (2025). Understanding One-Class Support Vector Machines. GeeksforGeeks. <https://www.geeksforgeeks.org/understanding-one-class-support-vector-machines/>
- [16] Scikit-learn. (2025). OneClassSVM. Scikit-learn. <https://scikit-learn.org/dev/modules/generated/sk-learn.svm.OneClassSVM.html>
- [17] Baeldung. (2025). One-Class SVM. Baeldung. <https://www.baeldung.com/cs/one-class-svm>
- [18] Medium. (2025). Anomaly detection using support vectors. Medium. <https://medium.com/@roshmitadey/anomaly-detect-ion-using-support-vectors-2c1b842213ed>
- [19] Anna Mrozek, Bartosz Panek, Przemysław Jura. (2025). Analiza zbioru pacjētów z SM chorych na COVID-19 modelem One-Class SVM. <https://github.com/przemyslaw-jura00/RozpoznawanieWzorcowGrupa4>



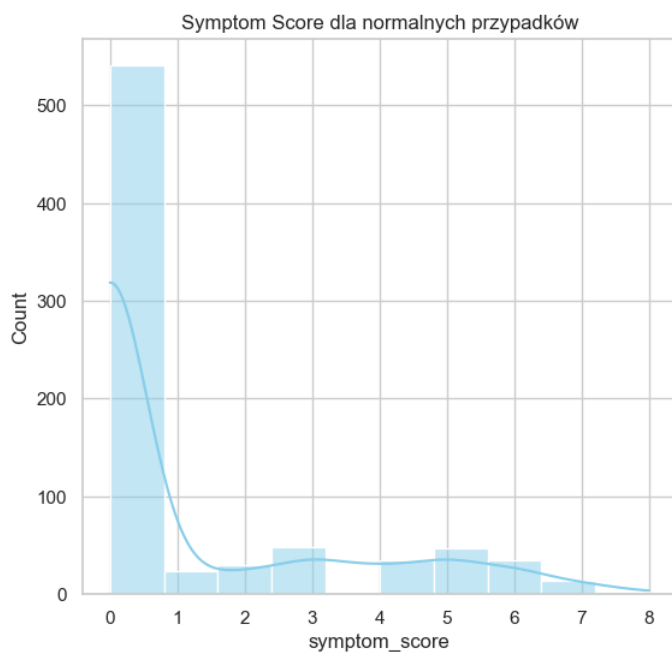
**Figure 3:** Wykryte anomalia vs Normalne próbki: -1: anomalia  
1: normalny przypadek



**Figure 4:** Rozkład ilości symptomów dla anomalii

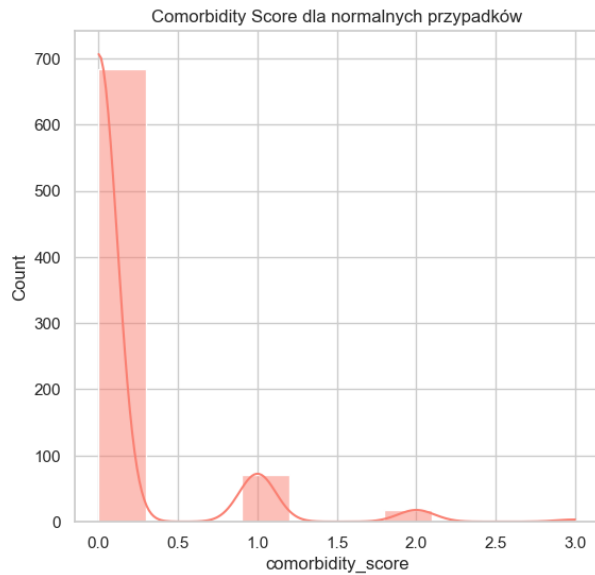


**Figure 5:** Rozkład ilości chorōb wspōistniejācych dla anomalii

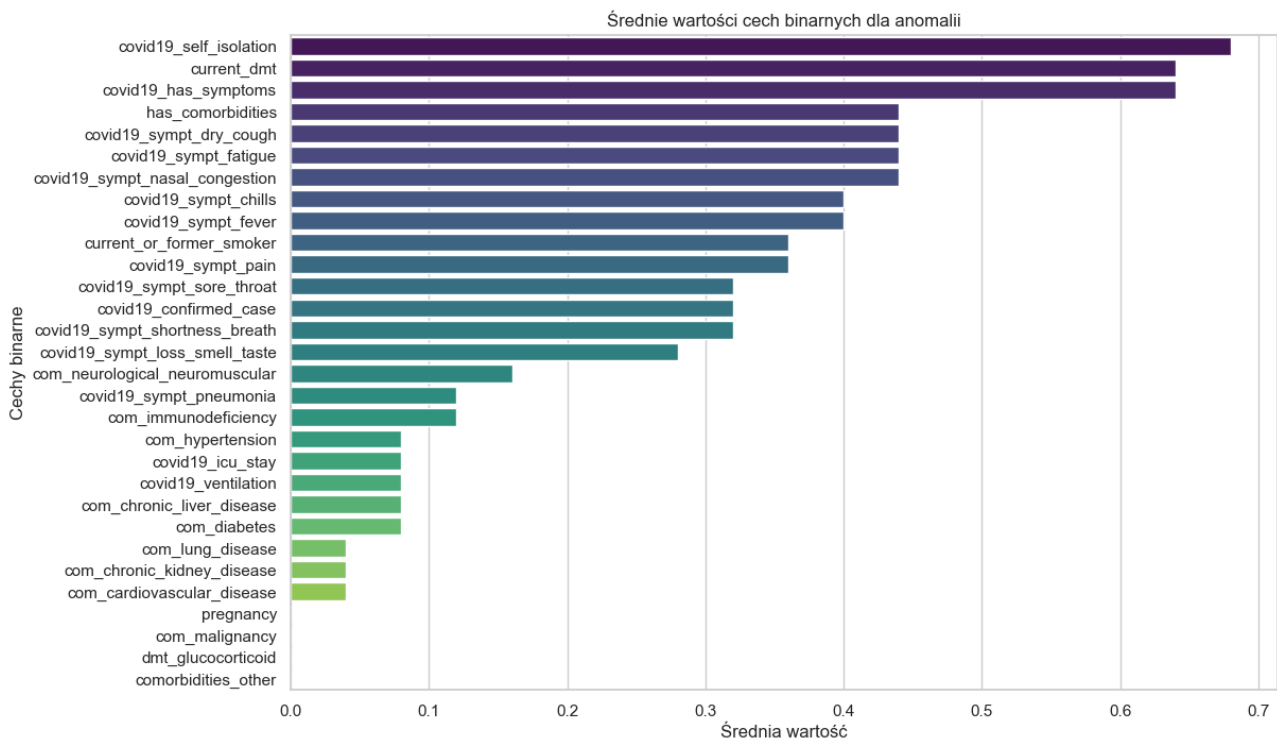


**Figure 6:** Rozkład ilości symptomōw dla przypadkōw normalnych

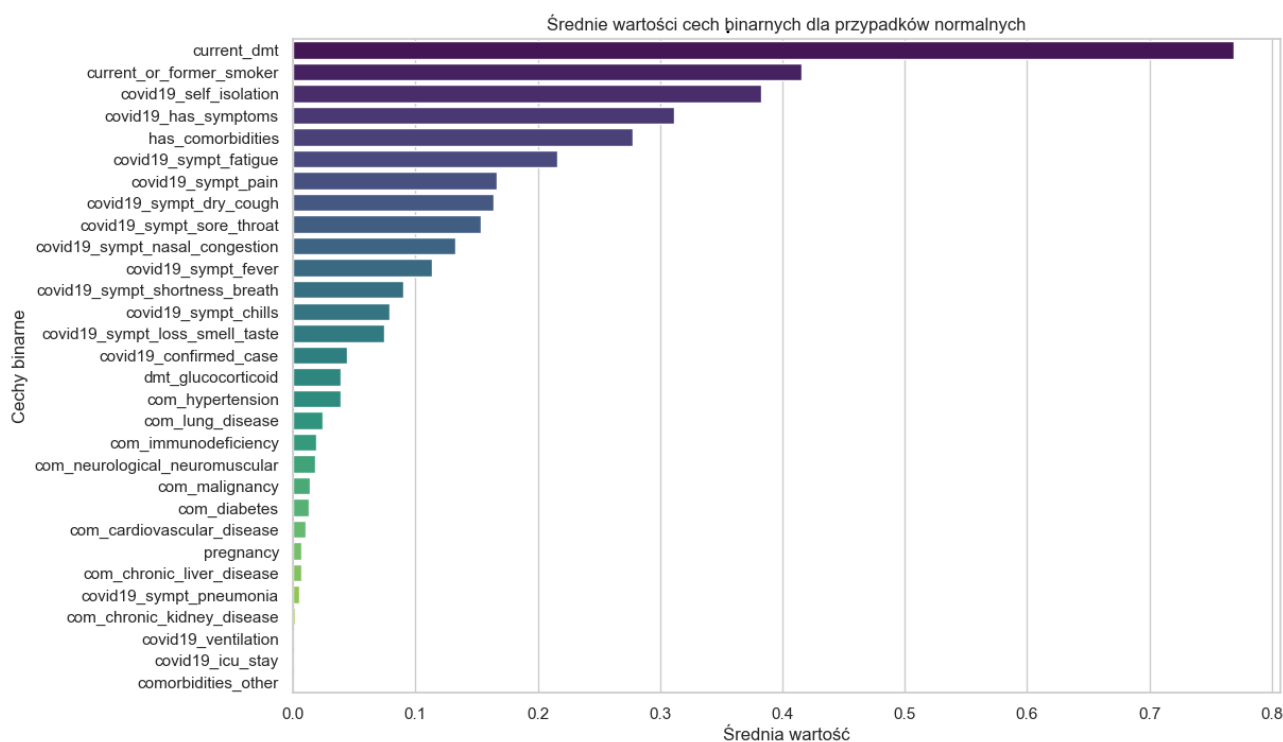




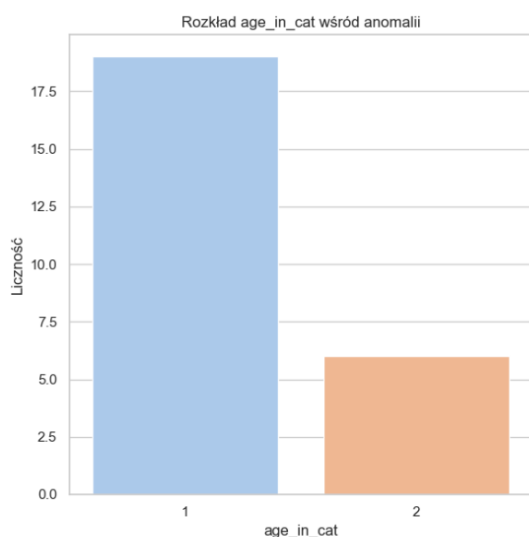
**Figure 7:** Rozkład ilości chorōb wspōistniejācych dla przypadkōw normalnych



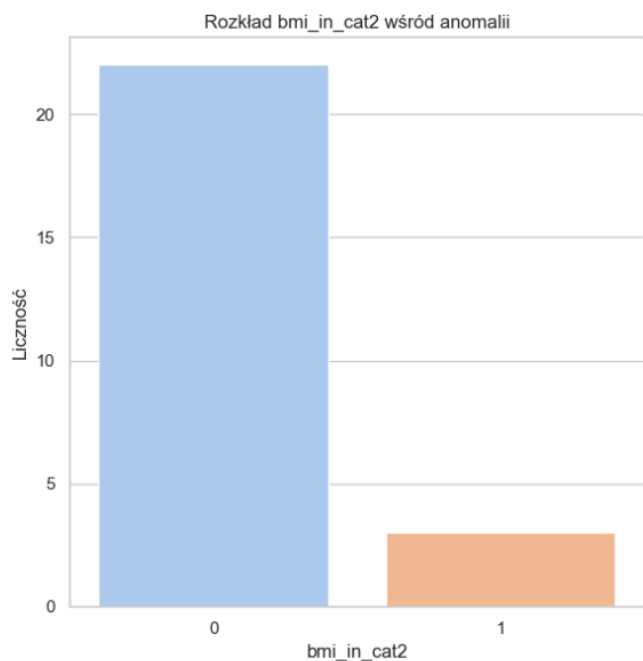
**Figure 8:** Średnie wartości charakterystycznych cech dla anomalii



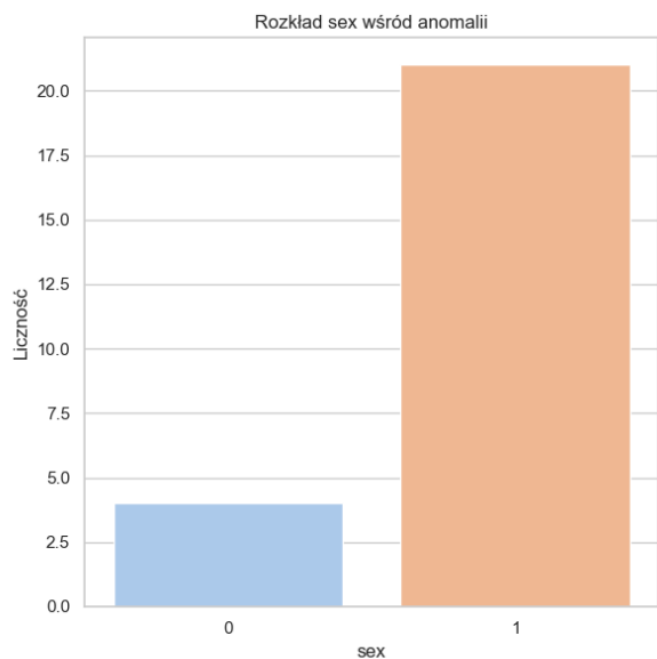
**Figure 9:** Średnie wartości charakterystycznych cech dla przypadkōw normalnych



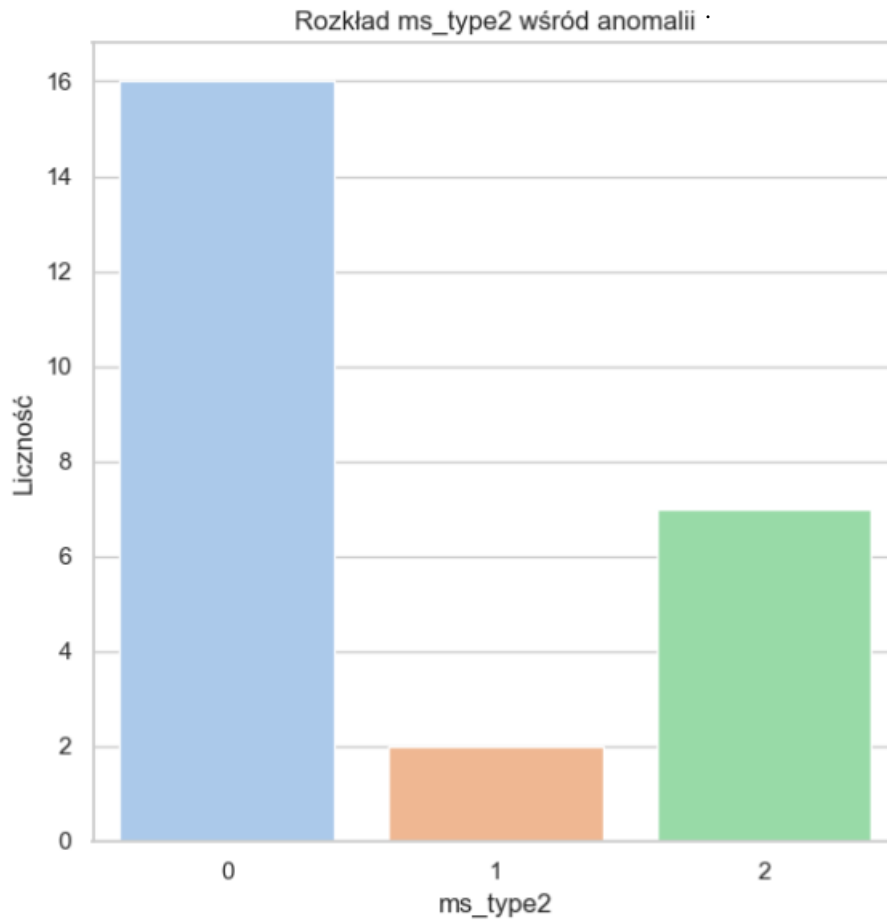
**Figure 10:** Rozkład wieku wśród anomalii: 0: jeśli zakres wieku mieści się w przedziale od 0 do <18. 1: jeżeli przedział wiekowy mieści się w przedziale od 18 do <=50 lat. 2: jeżeli przedział wiekowy mieści się w przedziale od 51 do <=70 lat. 3: jeśli przedział wiekowy wynosi 71 lat lub więcej..



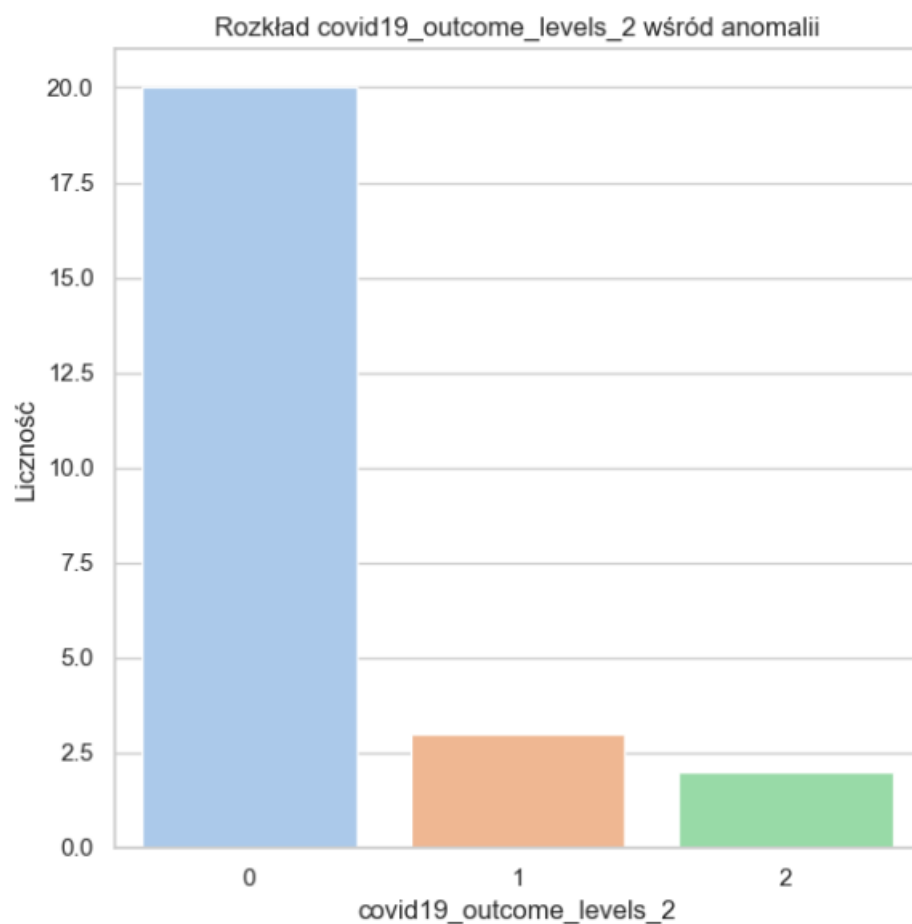
**Figure 11:** Rozkład bmi wśród anomalii: 0: not\_overweight: if BMI  $\leq 30$  kg/m<sup>2</sup> 1: overweight: if BMI  $> 30$  kg/m<sup>2</sup>.



**Figure 12:** Rozkład płci wśród anomalii: 0: mężczyźni 1: kobiety



**Figure 13:** Rozkład typu stwardnienia rozсіяnego wśród anomalii: 0: relapsing\_remitting: jeśli typ stwardnienia rozсіяnego to stwardnienie rozсіяne rzutowo-remisyjne (RRMS) 1: progresywny\_MS: jeśli typ stwardnienia rozсіяnego to wtórnie postępujące stwardnienie rozсіяne (SPMS) lub pierwotnie postępujące stwardnienie rozсіяne (PPMS) 2: inny: jeśli typ stwardnienia rozсіяnego to zespół izolowany klinicznie (CIS) lub pusty lub „niepewny”, w przypadku gdy pacjent lub lekarz nie był pewien.



**Figure 14:** Rozkłād hospitalizowanych przypadkōw: 0: Jeśli dana osoba ma Covid-19, ale nie była hospitalizowana. 1: Osoba ma Covid-19 i została hospitalizowana. 2: Osoba ma Covid-19, była hospitalizowana, przebywała na oddziale intensywnej terapii i/lub przebywała w ośrodku wentylacyjnym. 3: Osoba zmarła z powodu Covid-19 (nieobecna w tym zbiorze danych).