

Politechnika Wrocławska  
Wydział Podstawowych Problemów Techniki

# Jezyki programowania do zastosowań biomedycznych

KLASYFIKATOR WYPADKÓW

*Autorzy:*  
ZAKIA SHEFA #264476

Prowadzacy: dr hab. inż. Witold Dyrka

# Spis treści

<b>1</b>	<b>Wstęp</b>	<b>2</b>
<b>2</b>	<b>Zbiór danych</b>	<b>2</b>
<b>3</b>	<b>Wstępne przetwarzanie danych</b>	<b>3</b>
<b>4</b>	<b>Analiza danych</b>	<b>3</b>
4.1	Aktywności codzienne . . . . .	3
4.2	Upadki . . . . .	4
<b>5</b>	<b>Ekstrakcja cech</b>	<b>5</b>
5.1	Obliczanie wektorowych wielkości . . . . .	5
5.2	Okno czasowe i statystyki . . . . .	5
5.3	Zmienna celu . . . . .	6
<b>6</b>	<b>Podział na dane treningowe oraz testowe</b>	<b>6</b>
6.1	Podział losowy (Random Split) . . . . .	6
6.2	Podział na osoby (Split by Person) . . . . .	6
<b>7</b>	<b>Budowa Modelu</b>	<b>7</b>
7.1	K-Nearest Neighbors (KNN) . . . . .	7
7.2	Random Forest (RF) . . . . .	7
7.3	Support Vector Machines (SVM) . . . . .	7
<b>8</b>	<b>Eksperyment</b>	<b>8</b>
8.1	Podział losowy . . . . .	8
8.1.1	KNN . . . . .	8
8.1.2	RF . . . . .	8
8.1.3	SVM . . . . .	8
8.2	Podział na osoby . . . . .	9
8.2.1	KNN . . . . .	9
8.2.2	RF . . . . .	9
8.2.3	SVM . . . . .	9
<b>9</b>	<b>Podsumowanie</b>	<b>10</b>

# 1 Wstep

Celem pracy jest analiza, wybór oraz implementacja modelu klasyfikacji do wykrywania upadków u osób starszych z wykorzystaniem biosensorów wbudowanych w smartfony. Praca obejmuje przegląd metod uczenia maszynowego, przygotowanie zbioru danych na podstawie sygnałów z akcelerometru i żyroskopu, a następnie trenowanie modelu klasyfikacyjnego. Finalnym etapem jest integracja modelu z aplikacją mobilną działającą na systemie Android oraz testowanie skuteczności detekcji upadków.

## 2 Zbiór danych

Do realizacji projektu wykorzystano zbiór danych dostępny pod adresem: <https://www.kaggle.com/datasets/kmknation/mobifall-dataset-v20/data>.

Zbiór Mobifall zawiera dane sensoryczne pozyskane z akcelerometru, żyroskopu oraz czujnika orientacji. Dane zostały zebrane od 24 uczestników w wieku od 30 do 60 lat. Każda osoba wykonała zestaw określonych czynności, obejmujących zarówno codzienne aktywności, jak i symulowane upadki.

W zbiorze znajdują się także dodatkowe informacje demograficzne, takie jak: wzrost, waga, płeć oraz wiek badanych.

Zbiór jest podzielony na zestaw następujących aktywności:

### Codzienne aktywności (Activities of Daily Living)

Table 1: Zestaw codziennych aktywności

ID	Kod	Opis
1	STD	Stanie z delikatnymi ruchami
2	WAL	Normalne chodzenie
3	JOG	Bieganie (jogging)
4	JUM	Skakanie ciągle
5	STU	Wchodzenie po schodach (10 stopni)
6	STN	Schodzenie ze schodów (10 stopni)
7	SCH	Siadanie na krześle
8	CSI	Wsiadanie do samochodu
9	CSO	Wysiadanie z samochodu

### Upadki (Falls)

Table 2: Zestaw symulowanych upadków

ID	Kod	Opis
10	FOL	Upadek do przodu ze stania, z amortyzacją rękami
11	FKL	Upadek do przodu ze stania, pierwsze uderzenie kolanami
12	BSC	Upadek do tyłu podczas próby siadania na krześle
13	SDL	Upadek na bok ze stania, z ugięciem nóg

### 3 Wstępne przetwarzanie danych

Pierwszym krokiem w przetwarzaniu danych było ich odpowiednie przygotowanie do dalszej analizy. Ze względu na to, że dane zebrano w oddzielnych plikach dla każdego sensora, konieczne było ich połączenie w jeden spójny zbiór, co znacząco ułatwiło dalsze etapy pracy.

W trakcie integracji danych uwzględniono różnice w częstotliwości próbkowania pomiędzy akcelerometrem a żyroskopem. Akcelerometr (ACC) rejestrował dane z częstotliwością 50 Hz, natomiast żyroskop (GYRO) – z dużo wyższą, bo aż 2000 Hz. W związku z tym niektóre próbki żyroskopu zostały pominięte, a wartości z obu sensorów zostały połączone na podstawie najbliższych w czasie znaczników (timestampów).

Table 3: Częstotliwość próbkowania sensorów	
Sensor	Częstotliwość próbkowania
Akcelerometr (ACC)	0,02 s (20 ms) / 50 Hz
Żyroskop (GYRO)	0,0005 s (0,5 ms) / 2000 Hz

Na podstawie połączonych danych wygenerowano spójne pliki CSV z następującymi nagłówkami:

```
timestamp, x_acc, y_acc, z_acc, x_gyro, y_gyro, z_gyro, fall, activity, subject.id,
age, height, weight, gender
```

### 4 Analiza danych

W tej części przedstawiono przykładowe wykresy sygnałów sensorycznych (akcelerometru i żyroskopu), zarejestrowanych podczas wybranych aktywności codziennych oraz upadków. Dane przedstawiono w postaci wektorów wartości w funkcji czasu. Każdy wykres odpowiada jednemu pacjentowi. Poniżej znajdują się wizualizacje dla trzech różnych przypadków dla każdej aktywności i upadku.

#### 4.1 Aktywności codzienne

##### CSO – Car Step Out (Wysiadanie z samochodu)

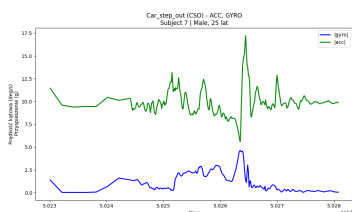


Figure 1: CSO – pacjent 7

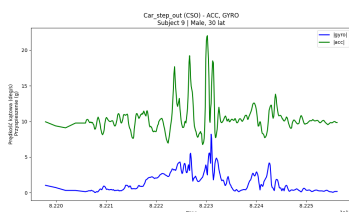


Figure 2: CSO – pacjent 9

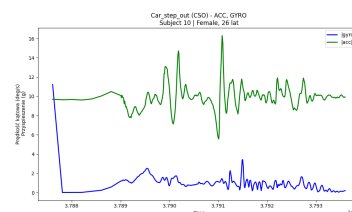


Figure 3: CSO – pacjent 10

##### JUM – Jumping (Skakanie)

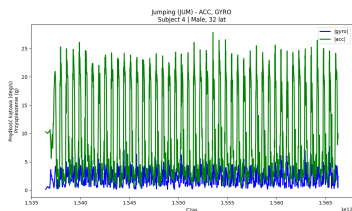


Figure 4: JUM – pacjent 4

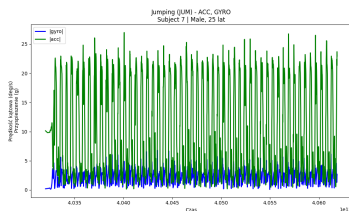


Figure 5: JUM – pacjent 7

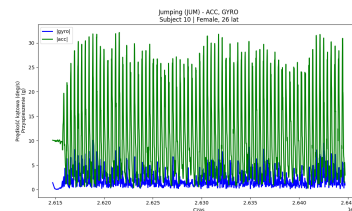


Figure 6: JUM – pacjent 10

## SCH – Sit Chair (Siadanie na krześle)

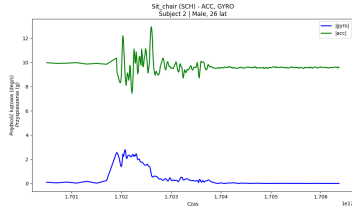


Figure 7: SCH – pacjent 2

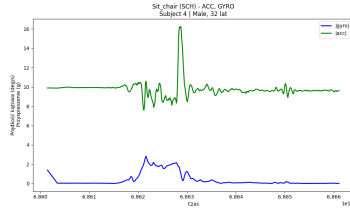


Figure 8: SCH – pacjent 4

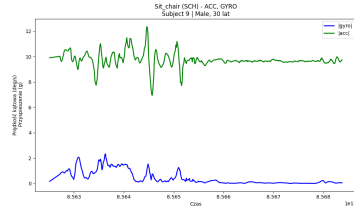


Figure 9: SCH – pacjent 9

Charakterystycznym w wynikach z akcelerometru wynikach jest to że nawet po największym ruchu sensora dane nie spadają szybko do stabilnego ruchu tylko zazwyczaj do mniejszego tak jak np na Wykresach SCH 7 8 9. Podobnie dane z gyrokopu są chaotyczne i widocznie pokazujące ciąg aktywności.

## 4.2 Upadki

### FOL – Forward Lying (Upadek do przodu z amortyzacją rękami)

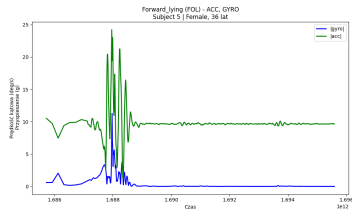


Figure 10: FOL – pacjent 5

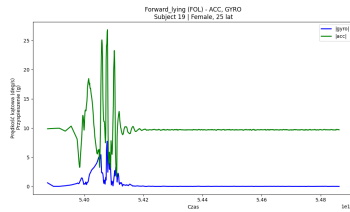


Figure 11: FOL – pacjent 19

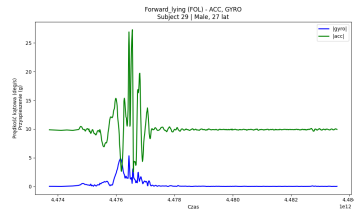


Figure 12: FOL – pacjent 29

### FKL – Front Knees Lying (Upadek do przodu z pierwszym uderzeniem kolan)

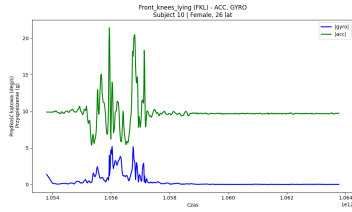


Figure 13: FKL – pacjent 10

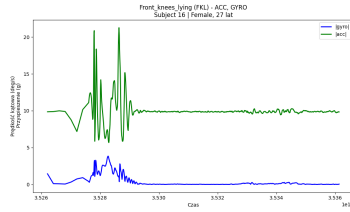


Figure 14: FKL – pacjent 16

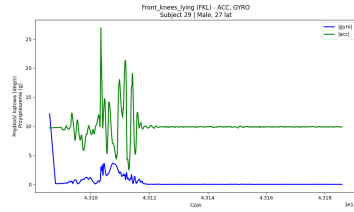


Figure 15: FKL – pacjent 29

### SDL – Sideward Lying (Upadek na bok)

Praktycznie wszystkie wyniki przedstawiają jeden duży skok a następnie 'wyciszenie'. podobnie w Przypadku gyrokopu zauważalne duże skoki danych a następnie wyrownanie

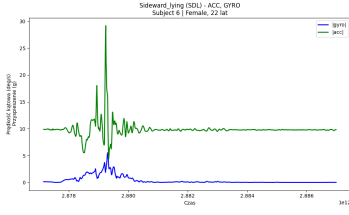


Figure 16: SDL – pacjent 6

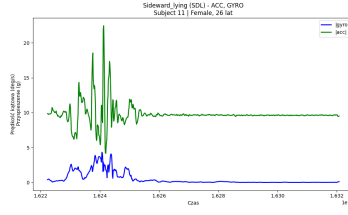


Figure 17: SDL – pacjent 11

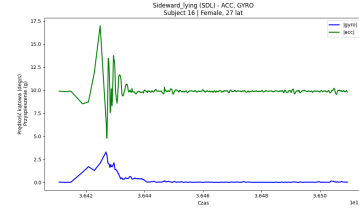


Figure 18: SDL – pacjent 16

## 5 Ekstrakcja cech

W celu umożliwienia skutecznej klasyfikacji zdarzeń, konieczne było przetworzenie surowych danych z czujników akcelerometru oraz żyroskopu na reprezentatywny zbiór cech numerycznych. Każda próbka zawierała pomiary przyspieszenia ( $x_{acc}$ ,  $y_{acc}$ ,  $z_{acc}$ ) oraz predkości katowej ( $x_{gyro}$ ,  $y_{gyro}$ ,  $z_{gyro}$ ).

### 5.1 Obliczanie wektorowych wielkości

Pierwszym krokiem było utworzenie pośrednich zmiennych reprezentujących wielkość wektorowa przyspieszenia oraz predkości katowej:

- **Moduł przyspieszenia ( $acc\_magnitude$ ):**

$$acc\_magnitude = \sqrt{x_{acc}^2 + y_{acc}^2 + z_{acc}^2}$$

Określa ogólna siłę działania sił inercyjnych i odzwierciedla intensywność ruchu.

- **Moduł predkości katowej ( $gyro\_magnitude$ ):**

$$gyro\_magnitude = \sqrt{x_{gyro}^2 + y_{gyro}^2 + z_{gyro}^2}$$

Umożliwia ocene rotacji urządzenia w przestrzeni 3D.

Następnie dla każdej z powyższych wielkości obliczono różnice pomiędzy kolejnymi próbkami, co pozwala uchwycić nagłe zmiany (cechy dynamiczne):

- $acc\_diff$  = różnica pomiędzy kolejnymi wartościami  $acc\_magnitude$ ,
- $gyro\_diff$  = różnica pomiędzy kolejnymi wartościami  $gyro\_magnitude$ .

### 5.2 Okno czasowe i statystyki

W celu ujednolicenia długości sekwencji czasowych oraz umożliwienia zastosowania klasyfikatorów klasycznych, dane zostały podzielone na nachodzące na siebie okna czasowe (tzw. *sliding window*):

- długość okna: 100 próbek (co odpowiada około 2 sekundom),
- przesunięcie okna: 10 próbek (aby uzyskać płynność analizy).

Dla każdej sekwencji w oknie obliczane są następujące cechy statystyczne:

- **Średnia wartość modułu przyspieszenia i predkości katowej ( $acc\_magnitude\_mean$ ,  $gyro\_magnitude\_mean$ )** – uśredniona siła ruchu w danym oknie,
- **Średnia różnic przyspieszenia i predkości katowej ( $acc\_diff\_mean$ ,  $gyro\_diff\_mean$ )** – zmienność sygnału,
- **Odchylenie standardowe modułów ( $acc\_magnitude\_std$ ,  $gyro\_magnitude\_std$ )** – miara rozrzutu wartości, odzwierciedlająca nieregularność ruchu.

Wprowadzono również atrybuty demograficzne:

- **age, height, weight** – wiek, wzrost i waga użytkownika (wartości maksymalne w oknie),
- **gender** – zakodowana liczbowo płeć (0 – mężczyzna, 1 – kobieta).

Cechy demograficzne te mogą mieć istotne znaczenie przy różnicowaniu wzorców ruchowych różnych użytkowników, np. ze względu na masę ciała lub sposób poruszania się.

### 5.3 Zmienna celu

Jako zmienna decyzyjna (*target*) wykorzystano binarną etykietę **fall**, wskazującą czy w danym oknie czasowym wystąpił upadek (wartość maksymalna tej kolumny w oknie czasowym – 1 oznacza co najmniej jeden upadek, 0 – brak).

## 6 Podział na dane treningowe oraz testowe

W celu oceny skuteczności działania klasyfikatora, przeprowadzono testy z wykorzystaniem dwóch różnych podejść do podziału zbioru danych: podziału losowego oraz podziału na podstawie identyfikatorów osób badanych. Oba podejścia pozwalają ocenić model w różnych kontekstach i weryfikują jego zdolność do generalizacji.

### 6.1 Podział losowy (Random Split)

W pierwszym podejściu zastosowano klasyczny podział zbioru danych na część treningową i testową w stosunku 70% do 30%. Dane były losowo dzielone przy użyciu funkcji `train_test_split()` z biblioteki `scikit-learn`:

```
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42)
```

Wybór stałej wartości `random_state=42` gwarantuje powtarzalność wyników. To podejście jest korzystne do szybkiej oceny ogólnej skuteczności klasyfikatora na przekrojowym zbiorze danych.

Należy jednak pamiętać, że dane nie są równomiernie rozłożone pomiędzy osoby – nie każdy badany dostarczył tyle samo próbek. W związku z tym, losowy podział może prowadzić do optymistycznych wyników, ponieważ przykłady z tej samej osoby mogą trafić zarówno do zbioru treningowego, jak i testowego.

### 6.2 Podział na osoby (Split by Person)

Drugie podejście opiera się na bardziej realistycznym scenariuszu: zakłada ono, że klasyfikator będzie używany przez osobę, która wcześniej nie występowała w danych treningowych. W tym celu dokonano podziału danych na podstawie identyfikatora osoby (`id`).

Zbiór testowy zawierał dane pochodzące od wybranej podgrupy osób (np. pierwsze 30% unikalnych `id`), a pozostałe osoby trafiły do zbioru treningowego:

```
unique_subjects = features_df['subject_id'].unique()

# Mieshamy osoby losowo
np.random.shuffle(unique_subjects)

# Wyznaczamy ile osób ma być w zbiorze testowym
n_testowych = int(len(unique_subjects) * procent_testowych)

test_subjects = unique_subjects[:n_testowych]
train_subjects = unique_subjects[n_testowych:]

# Podział danych na podstawie subject_id
train_df = features_df[features_df['subject_id'].isin(train_subjects)]
test_df = features_df[features_df['subject_id'].isin(test_subjects)]
```

Takie podejście jest znacznie bardziej wymagające dla modelu, ponieważ musi on uogólnić wiedzę zdobywaną na danych pewnych osób i zastosować ją do zupełnie nowych użytkowników. Testowanie z wykorzystaniem tego wariantu pozwala lepiej ocenić, jak model poradzi sobie w rzeczywistym zastosowaniu – np. jako część aplikacji mobilnej używanej przez nowych użytkowników.

## 7 Budowa Modelu

Do eksperymentów wybrano trzy algorytmy: **K-Nearest Neighbors (KNN)**, **Random Forest (RF)** oraz **SVC**. Wybór tych metod oparty był na ich skuteczności w problemach klasyfikacyjnych oraz zdolności do pracy z wielowymiarowymi danymi czasowymi.

### 7.1 K-Nearest Neighbors (KNN)

KNN to klasyfikator oparty na miarze odległości (np. euklidesowej) pomiędzy punktami w przestrzeni cech. Klasyfikacja nowej próbki odbywa się poprzez analizę  $k$  najbliższych sąsiadów w zbiorze treningowym — próbka zostaje przypisana do klasy, która dominuje wśród tych sąsiadów.

Główne cechy:

- Brak założeń o rozkładzie danych – metoda nieparametryczna.
- Naturalna obsługa danych wielocechowych – istotne dla zbioru danych zawierającego cechy czasowe i demograficzne.
- Wysoka interpretowalność – klasyfikacja oparta na intuicyjnym założeniu „bliskości”.

W eksperymentach stosowano różne wartości parametru  $k$  (liczby sąsiadów), z których najlepsze wyniki osiągnięto dla  $k = 3$  oraz  $k = 5$ . Dla poprawy jakości działania, dane zostały uprzednio standaryzowane (`StandardScaler` z biblioteki `scikit-learn`), co jest konieczne ze względu na wrażliwość KNN na różne skale cech.

### 7.2 Random Forest (RF)

Random Forest to model zespołowy (ensemble learning), który składa się z wielu drzew decyzyjnych tworzonych na losowych podzbiorach danych i cech. Końcowa predykcja modelu powstaje w wyniku głosowania większościowego.

Główne cechy:

- Wysoka odporność na przeuczenie dzięki losowości i agregacji wyników.
- Skuteczność w pracy z dużymi zbiorami danych i duża liczba cech.
- Możliwość interpretacji modelu poprzez analizę ważności cech.

Random Forest został wybrany jako bardziej zaawansowany klasyfikator do porównania z KNN, ze względu na swoją zdolność do modelowania nieliniowych zależności oraz odporność na szum w danych.

### 7.3 Support Vector Machines (SVM)

Algorytm **Support Vector Machines (SVC)** to jedna z najpopularniejszych metod klasyfikacyjnych, szczególnie skuteczna w przypadku danych o wysokiej wymiarowości. Jego głównym celem jest znalezienie hiperpowierzchni (tzw. hiperpłaszczyzny), która najlepiej rozdziela próbki należące do różnych klas, maksymalizując margines między nimi.

Główne cechy:

- Skuteczność w problemach klasyfikacyjnych, zwłaszcza przy ograniczonej liczbie próbek i dużej liczbie cech.
- Możliwość stosowania funkcji jadra (kernel), co pozwala na modelowanie nieliniowych granic decyzyjnych.
- Wysoka ogólność i elastyczność w dostosowywaniu do różnych typów danych.



## 8 Eksperyment

### 8.1 Podział losowy

Podział losowy polegał na odseparowaniu przypadków aktywności do zbioru testowego niezależnie od osób. Testy przeprowadzono dla długości okna 150 oraz przesunięcia okna o 50 próbek.

#### 8.1.1 KNN

	precision	recall	f1-score	support
0	0.96	0.98	0.97	4532
1	0.93	0.84	0.88	1264
accuracy			0.95	5796

Model KNN osiągnął bardzo dobre wyniki, zwłaszcza dla klasy codziennych aktywności (klasa 0), z wysoką precyzją i czułością. Dla klasy upadków (klasa 1), choć recall był nieco niższy (0.84), model skutecznie identyfikował większość przypadków.

#### 8.1.2 RF

	precision	recall	f1-score	support
0	0.98	0.98	0.98	4532
1	0.94	0.92	0.93	1264
accuracy			0.97	5796

Random Forest okazał się najlepszym modelem w tym scenariuszu, osiągając najwyższą dokładność (97%). Wysokie wartości precyzji i czułości dla obu klas wskazują na dużą skuteczność i równowagę w rozpoznawaniu zarówno upadków, jak i codziennych aktywności.

#### 8.1.3 SVM

	precision	recall	f1-score	support
0	0.93	0.97	0.95	4532
1	0.88	0.75	0.81	1264
accuracy			0.92	5796

Model SVM uzyskał nieco niższą skuteczność niż KNN i RF, zwłaszcza w wykrywaniu klasy 1, gdzie recall wyniósł 0.75. Choć dokładność ogólna nadal była wysoka (92%), SVM gorzej radził sobie z detekcją rzadkich przypadków upadków.

## 8.2 Podział na osoby

W tym podejściu osoby testowe zostały całkowicie wyłączone ze zbioru treningowego, co lepiej odwzorowuje realne zastosowanie modelu. Do testów wybrano 20% osób mających zarówno aktywności codzienne, jak i upadki.

### 8.2.1 KNN

	precision	recall	f1-score	support
0	0.64	0.24	0.35	1726
1	0.24	0.63	0.35	646
accuracy			0.35	2372

Powyższy wynik pokazuje, że model KNN rozpoznaje tylko niewielką część klas. Czulość (recall) dla upadków (klasa 1) jest względnie wysoka (0.63), ale precyzja bardzo niska (0.24), co oznacza dużą liczbę fałszywych alarmów.

### 8.2.2 RF

**Najlepszy przypadek:**

precision	recall	f1-score	support	
0	0.94	0.89	0.91	1618
1	0.75	0.85	0.79	629
accuracy			0.88	2247

W najbardziej korzystnym przypadku RF osiągnął bardzo dobre wyniki – dokładność 88% i wysokie wartości f1-score dla obu klas. To pokazuje, że model może skutecznie generalizować, nawet w warunkach testowania na nowych osobach.

**Gorszy przypadek:**

precision	recall	f1-score	support	
0	0.95	0.42	0.58	1694
1	0.39	0.94	0.55	654
accuracy			0.57	2348

W mniej korzystnym scenariuszu model nadal dobrze wykrywał upadki (recall = 0.94), ale miał bardzo niską skuteczność w rozpoznawaniu aktywności codziennych, co prowadziło do wielu fałszywych alarmów.

### 8.2.3 SVM

W celu poprawy detekcji rzadkiej klasy upadków, zastosowano zwiększoną wagę tej klasy podczas trenowania modelu.

**Wagi klasy 1: 5**

precision	recall	f1-score	support	
0	0.81	0.45	0.58	1726
1	0.33	0.72	0.45	646
accuracy			0.53	2372

Model lepiej identyfikował upadki (recall = 0.72), jednak kosztem licznych fałszywych alarmów (niska precyzja 0.33). Skuteczność dla klasy 0 również pozostała niezadowolająca.

**Wagi klasy 1: 10**

precision	recall	f1-score	support	
0	0.91	0.63	0.74	1694
1	0.47	0.85	0.60	654
accuracy			0.69	2348

Zwiększenie wagi klasy upadków do 10 poprawiło wyniki modelu – recall dla upadków wzrósł do 0.85, a ogólna dokładność osiągnęła 69%. Mimo to, precyzja nadal pozostaje umiarkowana (0.47), co może prowadzić do fałszywych powiadomień.

## 9 Podsumowanie

Podejście z podziałem na osoby znacząco utrudniało zadanie klasyfikatorom – pokazując rzeczywiste możliwości generalizacji. KNN całkowicie zawiódł w tym scenariuszu. SVM wykazał poprawę po odpowiednim dostrojeniu wag, ale nadal nie dorównał skuteczności RF. Random Forest, mimo zmienności w wynikach, wciąż osiągał najlepszy kompromis między precyzją a czułością. Dla podziału losowego wszystkie algorytmy osiągnęły bardzo wysoką dokładność (powyżej 90%). Klasa upadków została dobrze wykryta, jednak wyniki mogą być zawyżone z powodu obecności danych tych samych osób w zbiorach treningowym i testowym. Może to prowadzić do przeuczenia na indywidualne cechy (np. wzrost, sposób chodzenia), co w praktycznym zastosowaniu może obniżać skuteczność.

Na podstawie przeprowadzonych eksperymentów, najlepsze rezultaty osiągnął model **Random Forest** — zarówno w scenariuszu podziału losowego, jak i realistyczniejszym podziale na osoby. RF okazał się najstabilniejszy, a także najbardziej zrównoważony w wykrywaniu upadków i unikaniu fałszywych alarmów.

Biorąc pod uwagę końcowe zastosowanie modelu w aplikacji wspierającej bezpieczeństwo osób starszych, kluczowe jest skuteczne wykrywanie upadków, nawet kosztem umiarkowanej liczby fałszywych alarmów. Random Forest, ze względu na swoją elastyczność i wysoką dokładność, jest najbardziej obiecującym kandydatem do implementacji.