

# Preprocessing danych – Titanic dataset

1. Wstęp.....	3
2. Teoria – podstawy teoretyczne.....	4
3. Preprocessing .....	6
4. Analiza - część analityczna .....	8
5. Wyniki.....	10
6. Podsumowanie – wnioski .....	12
7. Bibliografia.....	15
8. Kod.....	16

# 1. Wstęp

Celem niniejszego sprawozdania jest przeprowadzenie kompletnego procesu wstępniego przetwarzania (preprocessingu) oraz eksploracyjnej analizy danych (EDA) słynnego zbioru Titanic (plik titanic\_new.csv). Przygotowanie danych ma kluczowe znaczenie, ponieważ poprawia ich jakość i pozwala na wiarygodne modelowanie — na przykład budowę klasyfikatora przewidującego przeżycie pasażerów.

Zakres pracy obejmuje:

1. **Identyfikację, diagnostykę i usunięcie problemów w danych** – braki, wartości odstające, niespójności.
2. **Transformację i kodowanie zmiennych** oraz, gdy zachodzi taka potrzeba, **redukcję wymiarów** w celu uproszczenia struktury danych.
3. **Wyznaczenie i interpretację kluczowych statystyk opisowych** wraz z przejrzystymi wizualizacjami rozkładów i relacji między cechami.
4. **Sformułowanie wniosków i rekomendacji** dotyczących dalszych analiz i modelowania predykcyjnego.

Tak uporządkowana sekwencja działań pozwoli nie tylko poprawić jakość samych danych, lecz także odsłonić ukryte zależności i zapewnić solidne fundamenty pod kolejne etapy pracy analitycznej.

## 2. Teoria – podstawy teoretyczne

Wstępne przetwarzanie (preprocessing) to zestaw technik, których zadaniem jest zamiana surowego, często niespójnego zbioru w spójny, kompletny i analitycznie użyteczny materiał. Zgodnie ze slajdami 9\_Preprocessing.pdf wyróżniamy cztery główne kategorie działań, z których trzy zostały zastosowane w niniejszym projekcie (ostatnia – integracja – nie była potrzebna, ponieważ pracujemy na pojedynczym pliku titanic\_new.csv).

Kategoria	Kluczowe operacje	Cel
<b>Czyszczenie danych</b>	<ul style="list-style-type: none"> <li>• usuwanie duplikatów</li> <li>• uzupełnianie braków (lub usuwanie rekordów)</li> <li>• detekcja i korekta szumu oraz anomalii</li> <li>• poprawa lub konwersja typów danych</li> </ul>	Usunięcie błędów i niespójności, by uniknąć zafałszowania wyników analizy.
<b>Transformacja</b>	<ul style="list-style-type: none"> <li>• identyfikacja i obsługa wartości odstających</li> <li>• kodowanie zmiennych kategorycznych</li> <li>• skalowanie / standaryzacja</li> <li>• rozdzielenie złożonych kolumn na atomowe składniki</li> </ul>	Ujednolicenie reprezentacji cech oraz dostosowanie ich do wymagań algorytmów.
<b>Redukcja wymiarów</b>	<ul style="list-style-type: none"> <li>• eliminacja nieistotnych lub silnie skorelowanych kolumn</li> <li>• łączenie rzadkich kategorii</li> <li>• ewentualnie metody takie jak PCA</li> </ul>	Zmniejszenie liczby cech w celu uproszczenia modelu, ograniczenia przeuczenia i przyspieszenia obliczeń.
<b>Integracja (pominięta)</b>	<ul style="list-style-type: none"> <li>• scalanie danych z wielu źródeł</li> </ul>	Zbudowanie jednolitego repozytorium informacji, gdy dane pochodzą z różnych plików lub baz.

## Dlaczego preprocessing jest niezbędny?

- Poprawia **wiarygodność i jakość** danych, co przekłada się na rzetelniejsze statystyki i modele predykcyjne.
- Ogranicza **wpływ błędów pomiaru i szumu**, które mogłyby zawyżyć lub zaniżyć wyniki.
- Redukuje **złożoność** zbioru, co przyspiesza trening modeli i ułatwia interpretację rezultatów.

W projekcie wszystkie opisane etapy zrealizowano w języku **R** na bazie pojedynczego źródła (titanic\_new.csv). Tak przygotowany, oczyszczony i zoptymalizowany zestaw stanowi solidny fundament do dalszych prac analitycznych, takich jak eksploracja zmiennych i budowa klasyfikatora przewidującego przeżycie pasażera.

### 3. Preprocessing

Poniżej przedstawiono kompletny, uporządkowany pipeline przygotowania danych Titanic (plik titanic\_new.csv) do dalszej analizy i modelowania. Każdy krok zawiera krótkie uzasadnienie, dzięki czemu łatwo prześledzić, dlaczego wykonano daną operację i jak poprawia ona jakość zbioru danych.

Krok	Co zrobiono	Dlaczego
1. Usunięcie kolumn nieistotnych	Skasowano PassengerId, Name, Ticket, Cabin	Unikalne lub rzadkie ciągi znaków nie wspierają predykcji przeżycia; redukcja wymiaru ogranicza szum i przyspiesza wizualizacje.
2. Konwersja typów na kategorie	Survived, Pclass, Sex, Embarked → <i>factor</i>	Kategoryczne przechowywanie zmiennych obniża zużycie pamięci i ułatwia tworzenie wykresów oraz modelowanie.
3. Obsługa braków danych	- Age uzupełniono medianą (28 lat) - Embarked uzupełniono dominującą wartością S	Imputacja zachowuje liczebność próby i zapobiega odrzuceniu informacyjnych rekordów.
4. Kontrola duplikatów	Sprawdzono pełną identyczność rekordów – <b>duplikatów brak</b>	Zapobiega sztucznemu zawyżaniu częstości obserwacji.
5. Detekcja wartości odstających	Dla Age i Fare przyjęto do 99. percentyla	Stabilizuje algorytmy wrażliwe na ekstrema, zmniejsza wariancję modeli.
6. Kodowanie zmiennych kategorycznych	- Sex: female = 1, male = 0 - Embarked: kodowanie one-hot (S, C, Q)	Większość algorytmów ML wymaga reprezentacji numerycznej; one-hot zapobiega narzucaniu sztucznego porządku.
7. Skalowanie cech numerycznych	Standaryzacja Age, Fare, SibSp, Parch ( $\mu = 0, \sigma = 1$ )	Metody oparte na odległościach (np. K-NN, SVM, PCA) działają poprawniej na danych o jednolitych skalach.
8. Redukcja rzadkich kategorii (opcjonalnie)	Łączenie bardzo rzadkich wartości w kategoriach (gdy < 1 % obserwacji)	Zapewnia stabilniejsze oszacowania i usprawnia wizualizacje.

Efektem powyższych czynności jest oczyszczony i zoptymalizowany zbiór **df\_clean** o wymiarach  $891 \times 8$ , gotowy do:

- eksploracyjnej analizy zmiennych i zależności,
- budowy modeli predykcyjnych (np. klasyfikatora przeżycia pasażera),
- dalszych eksperymentów, takich jak inżynieria cech czy walidacja krzyżowa.

Dzięki konsekwentnemu zastosowaniu opisanych technik preprocessingowych zbiór danych charakteryzuje się **spójnością, kompletnością i mniejszą złożonością**, co stanowi solidny fundament dla kolejnych etapów pracy analitycznej.

# 4. Analiza - część analityczna

## Analiza opisowa i eksploracyjna

Po zakończeniu preprocessingu na zbiorze **df\_clean (891 × 8)** przeprowadzono wstępne analizy statystyczne i wizualizacje, aby lepiej zrozumieć strukturę danych i kluczowe czynniki wpływające na przeżycie pasażerów.

Obszar	Główne wyniki	Krótką interpretacja
<b>Rozkład zmiennej docelowej</b>	38 % pasażerów przeżyło • kobiety: <b>74 %</b> • mężczyźni: <b>19 %</b>	Potwierdza regułę „women and children first” oraz sugeruje, że płeć jest silnym predyktorem przeżycia.
<b>Statystyki opisowe</b>	• Mediana wieku: <b>28 lat</b> (najwięcej osób 20–30 l.) • <i>Fare</i> : skośny rozkład, 75 % wartości < 50, nieliczne obserwacje > 200	Sugeruje potrzebę log-transformacji lub przycięcia <i>Fare</i> w modelach wrażliwych na skośność.
<b>Korelacje z Survived</b>	• Najsilniejsze dodatnie: wysoka cena biletu ( <i>Fare</i> ), klasa 1 ( <i>Pclass</i> = 1) • Najsilniejszy ujemny: płeć męska ( <i>Sex</i> = 0)	Wysoka cena i klasa wskazują na lepszy dostęp do łodzi ratunkowych; płeć męska obniża szanse przeżycia.
<b>Wizualizacje</b>	• Histogram wieku – szczyt w przedziale 20–30 l. • Boxplot <i>Fare</i> vs <i>Survived</i> – wyższe ceny u ocalałych • Mapa ciepła korelacji – potwierdza powyższe zależności	Wizualizacje odsłaniają wzorce trudne do zauważenia w tabelach i wspierają wybór cech do modelu.
<b>Przeżycie wg płyci i klasy</b>	• Klasa 1: <b>63 %</b> ocalałych • Klasa 3: <b>24 %</b> ocalałych	Kombinacja niskiej klasy i płci męskiej dramatycznie obniża szanse przeżycia.

Wnioski operacyjne:

1. **Płeć, klasa podróży i cena biletu** to najsilniejsze pojedyncze predyktory; warto je zachować w modelu bez transformacji kategorii.
2. Skośność zmiennej *Fare* i obecność odstających wartości sugerują **transformację logarytmiczną** lub dalsze przycięcie przy budowie niektórych modeli.
3. Relatywnie wysoka korelacja Pclass–Fare podpowiada, by sprawdzić **wielokoliniowość**, zwłaszcza w modelach liniowych.
4. Dalszy kierunek pracy: implementacja i porównanie modeli (regresja logistyczna, drzewa decyzyjne, XGBoost, sieci neuronowe) z **walidacją krzyżową** i miarami takimi jak AUC-ROC.

Tak skondensowana analiza dostarcza nie tylko liczbowych podsumowań, ale i praktycznych wskazówek, na których cechach skupić się w kolejnych etapach budowy predyktora przeżycia pasażerów Titanica.

## 5. Wyniki

Wyniki po pełnym preprocessingu zbioru *Titanic* otrzymano **891 rekordów i 8 zmiennych**, całkowicie pozabawionych braków i gotowych do modelowania.

Obszar	Najważniejsze wyniki	Znaczenie praktyczne
<b>Braki danych</b>	<ul style="list-style-type: none"> <li>Imputacja wieku (medianą = 28 lat) obniżała odsetek wartości brakujących z 19 % do 0 %.</li> </ul>	Zachowano wszystkie obserwacje bez wprowadzania biasu związanego z usuwaniem wierszy.
<b>Standaryzacja</b>	<ul style="list-style-type: none"> <li>Po skalowaniu średnia kluczowych zmiennych numerycznych <math>\approx 0</math>, a odchylenie standardowe = 1.</li> </ul>	Dane są przygotowane do algorytmów opartych na odległości (SVM, K-NN, PCA).
<b>Wartości odstające</b>	<ul style="list-style-type: none"> <li>Przycięcie <i>Fare</i> do 99. percentyla ograniczyło maksymalną opłatę z 512 £ do 263 £.</li> </ul>	Zmniejszono wpływ ekstremów na parametry modeli wrażliwych na wariancję.
<b>Statystyki opisowe</b>	<ul style="list-style-type: none"> <li>Medianą wieku: <b>28 lat</b></li> <li>Ogólny wskaźnik przeżycia: <b>38,1 %</b></li> </ul>	Potwierdza młody profil wiekowy większości pasażerów i relatywnie niski odsetek ocalonych.
<b>Przeżywalność wg płci</b>	<ul style="list-style-type: none"> <li>Kobiety: <b>74,0 %</b></li> <li>Mężczyźni: <b>18,6 %</b></li> </ul>	Test $\chi^2$ : $p < 0,001 \rightarrow$ płeć to najsilniejszy pojedynczy czynnik przeżycia.
<b>Przeżywalność wg klasy</b>	<ul style="list-style-type: none"> <li>I klasa: 62,6 % •</li> <li>II klasa: 47,3 % •</li> <li>III klasa: 23,9 %</li> </ul>	Wyższa klasa i wyższa cena biletu ( <i>Fare</i> ) istotnie zwiększały szanse przeżycia.
<b>Korelacje z Survived</b>	<ul style="list-style-type: none"> <li>Najsilniejsze: płeć (ujemna korelacja dla mężczyzn), klasa 1, wysoki <i>Fare</i>.</li> </ul>	Kierunkuje wybór zmiennych w modelach predykcyjnych.

Kluczowe wnioski:

1. **Płeć** pozostaje najważniejszym predyktorem przeżycia, a istotność statystyczna ( $p < 0,001$ ) czyni ją obowiązkową w każdym modelu.
2. **Klasa podróży i cena biletu** odzwierciedlają status ekonomiczny, który również silnie koreluje z przeżyciem.
3. Rozkłady **Age** i **Fare** są asymetryczne; przy dalszym modelowaniu warto rozważyć transformację logarytmiczną *Fare* oraz ew. segmentację wieku.
4. Kompleksowa redukcja wymiaru i eliminacja odstających wartości poprawiły stabilność danych, co powinno przełożyć się na **lepszą wydajność i interpretowalność modeli** (regresja logistyczna, drzewa, XGBoost).

Zestawione wyniki dostarczają solidnych podstaw do kolejnego etapu – budowy i oceny modeli predykcyjnych prognozujących przeżycie pasażerów Titanica.

# 6. Podsumowanie – wnioski

## Podsumowanie, możliwe analizy i niezbędne dane

Co osiągnięto?

Starannie przeprowadzony preprocessing – obejmujący czyszczenie, imputację, usunięcie odstających wartości i standaryzację – przekształcił zbiór Titanic w spójny, kompletny i statystycznie stabilny zestaw  $891 \times 8$ . Dzięki temu analiza jest wiarygodna, a wnioski nieobciążone błędami danych.

Najważniejsze predyktory przeżycia:

- **Płeć (Sex)** – najsilniejszy czynnik ( $\chi^2$ ,  $p < 0,001$ ).
- **Klasa podróży (Pclass)** – wyższa klasa  $\Rightarrow$  większa szansa przeżycia.
- **Cena biletu (Fare)** – dodatnia korelacja z przeżyciem (odzwierciedla status ekonomiczny).
- **Wiek (Age)** – umiarkowany wpływ, szczególnie przy analizie interakcji z płcią.
- **Port zaokrętowania (Embarked)** – słabszy, ale przydatny w modelach z wieloma cechami. Uzupełniająco warto tworzyć zmienne pochodne, np. rodzina na pokładzie (FamilyOnBoard = SibSp + Parch > 0).

Jaką analizę można teraz wykonać?

Typ analizy	Cel	Kluczowe zmienne
Model klasifykacyjny (np. regresja logistyczna, drzewa, XGBoost)	Predykcja prawdopodobieństwa przeżycia każdego pasażera	Sex, Pclass, Fare, Age, Embarked, SibSp, Parch, zmienne pochodne
Walidacja krzyżowa + strojenie hiperparametrów	Porównanie skuteczności różnych algorytmów i ustawień	Te same zmienne co wyżej; metryki (AUC-ROC, dokładność, recall)

Analiza interakcji (np. współczynnik w regresji logistycznej, SHAP w XGBoost)	Zrozumienie, jak wpływ dwóch cech naraz (np. Sex × Age) zmienia szanse przeżycia	Sex, Age, Pclass, Fare
Segmentacja (klasteryzacja K-means / hierarchiczna)	Identyfikacja homologicznych grup pasażerów o podobnym profilu ryzyka	Zmienne znormalizowane: Age, Fare, Pclass, FamilyOnBoard
Analiza przetrwania (Survival Analysis, np. model Coxa)	Oszacowanie „czasu do zdarzenia”, czyli prawdopodobieństwa przeżycia w kolejnych minutach katastrofy	Wymaga przybliżonego czasu ewakuacji (jeśli dostępny) + Sex, Pclass, Age, itd.

Dlaczego akurat te dane?

- Predykcyjne – Sex, Pclass, Fare, Age i Embarked mają największą moc rozdzielczą wg testów statystycznych i korelacji.
- Kompletne – braków już nie ma; imputacja wieku i dominanta dla Embarked przywróciły 100 % rekordów.
- Zestandardyzowane – standaryzacja zmiennych numerycznych eliminuje problemy związane ze skalą przy algorytmach opartych na odległości.
- Oczyszczone ze skrajności – przycięcie Fare do 99. percentyla oraz kontrola wieku redukują ryzyko, że pojedyncze obserwacje zdominują proces uczenia.

Rekomendowane kolejne kroki

- Zbudować i porównać co najmniej trzy modele (regresja logistyczna, XGBoost, las losowy) z walidacją 10-krotną.
- Dodać cechy inżynieryjne (np. Title z nazwisk lub FamilyOnBoard) i sprawdzić, czy poprawiają AUC-ROC.

- Zbadać interakcje – szczególnie Sex × Age oraz Pclass × Fare.
- Przeprowadzić interpretację modelu (SHAP, permutacja ważności) w celu wyjaśnienia wpływu poszczególnych cech.
- Rozszerzyć imputację o metody wielokrotnego uzupełniania (MICE), aby uwzględnić niepewność imputowanych wartości.

Dzięki tak przygotowanemu zbiorowi danych i jasnej mapie dalszych analiz można zbudować dokładny, przejrzysty i interpretowalny model przewidujący przeżycie pasażerów Titanica — lub poszerzyć badanie o nowe, bardziej złożone pytania badawcze.

## 7. Bibliografia

1. **Maciejko M.** Preprocessing danych – slajdy wykładowe (9\_Preprocessing.pdf).
2. **Kaggle** – Titanic Dataset, oprac. S. Seaborn, 2019.
3. **Analytics Vidhya** – A Comprehensive Guide to Data Pre-processing

## 8. Kod

#Titanic Preprocessing and EDA Pipeline in R

#Wczytanie bibliotek

```
library(tidyverse) library(corrplot)
```

#1. Wczytanie danych

Zakładamy, że plik titanic\_new.csv znajduje się w katalogu roboczym

```
df<- read.csv("titanic_new.csv", stringsAsFactors = FALSE)
```

#2. Usunięcie kolumn nieistotnych

PassengerId, Name, Ticket, Cabin zawierają unikalne lub zbędne informacje

```
df <- df %>% select(-PassengerId, -Name, -Ticket, -Cabin)
```

#3. Konwersja typów na faktory

```
df <- df %>% mutate( Survived = factor(Survived, levels = c(0,1), labels = c("No","Yes")), Pclass = factor(Pclass, levels = c(1,2,3)), Sex = factor(Sex), Embarked = factor(Embarked) )
```

#4. Imputacja braków danych

#4.1. Wiek: mediana

```
i <- median(df$Age, na.rm = TRUE) df$Age[is.na(df$Age)] <- i
```

#4.2. Port zaokrągiania: dominanta

```
mode_emb <- df %>% filter(!is.na(Embarked)) %>% count(Embarked) %>%
arrange(desc(n)) %>% slice(1) %>% pull(Embarked)
df$Embarked[is.na(df$Embarked)] <- mode_emb
```

## #5. Kontrola duplikatów

```
Usunięcie identycznych rekordów. df <- df %>% distinct()
```

## #6. Detekcja i przycięcie wartości odstających do 99. percentyla

```
caps <- df %>% summarize( age_cap = quantile(Age, 0.99), fare_cap =
quantile(Fare, 0.99) ) df <- df %>% mutate( Age = ifelse(Age > caps$age_cap,
caps$age_cap, Age), Fare = ifelse(Fare > caps$fare_cap, caps$fare_cap, Fare) )
```

## #7. Kodowanie zmiennych kategorycznych

### #7.1. Sex: female = 1, male = 0

```
df <- df %>% mutate(Sex = ifelse(Sex == "female", 1, 0))
```

### #7.2. Embarked: one-hot encoding

```
dummies <- model.matrix(~ Embarked - 1, data = df) %>% as.data.frame() df
<- bind_cols(df %>% select(-Embarked), dummies)
```

## #8. Skalowanie cech numerycznych

```
num_vars <- c("Age", "Fare", "SibSp", "Parch") df[num_vars] <-
scale(df[num_vars])
```

## #Oczyszczony zbiór danych

```
df_clean <- df
```

## #Zapis oczyszczonego zbioru do pliku (opcjonalnie)

```
write.csv(df_clean, "titanic_clean.csv", row.names = FALSE)
```

----- Eksploracyjna Analiza Danych (EDA) -----

#Rozkład zmiennej docelowej (przeżycie)

```
surv_dist <- prop.table(table(df_clean$Survived)) * 100 print(surv_dist)
```

#Przeżywalność wg płci

```
tab_sex <- prop.table(table(df_clean$Sex, df_clean$Survived), 1) * 100  
print(tab_sex)
```

#Przeżywalność wg klasy podróży

```
tab_class <- prop.table(table(df_clean$Pclass, df_clean$Survived), 1) * 100  
print(tab_class)
```

#Statystyki opisowe dla Age i Fare

```
print(summary(df_clean$Age)) print(summary(df_clean$Fare))
```

#Korelacje zmiennych z Survived

#Zamiana Survived na 0/1 dla korelacji

```
corr_data <- df_clean %>% mutate(Survived = ifelse(Survived == "Yes", 1, 0))  
%>% select(Survived, all_of(num_vars), starts_with("Embarked")) corr_matrix  
<- cor(corr_data) corrplot(corr_matrix, method = "color", tl.cex = 0.8,  
addCoef.col = "black")
```

#Wizualizacje

#Histogram wieku

```
ggplot(df_clean, aes(x = Age)) + geom_histogram(binwidth = 5, fill = "lightblue", color = "black") + ggtitle("Histogram wieku po skalowaniu") + theme_minimal()
```

```
#Boxplot ceny biletu vs przeżycie
```

```
ggplot(df_clean, aes(x = Survived, y = Fare)) + geom_boxplot() + ggtitle("Boxplot ceny biletu wg przeżycia") + theme_minimal()
```

```
# Mapa ciepła korelacji
```

```
corrplot(corr_matrix, method = "heatmap", tl.cex = 0.8)
```

```
#Wykres przeżywalności wg płci i klasy
```

```
ggplot(df_clean, aes(x = Pclass, fill = Survived)) + geom_bar(position = "fill") + facet_wrap(~ Sex, labeller = labeller(Sex = c(0 = "Male", 1 = "Female")))) + ylab("Proporcja") + ggtitle("Przeżywalność wg klasy i płci") + theme_minimal()
```