

# Eksploracja danych - etap 1

Krzysztof Nasuta 193328, Filip Dawidowski 193433, Aleks Iwicki 193354

## 1. Ogólny opis zbioru

Departament Transportu Stanów Zjednoczonych (DOT), za pośrednictwem Biura Statystyki Transportu, monitoruje punktualność krajowych lotów realizowanych przez dużych przewoźników lotniczych. Zbiorcze informacje na temat liczby lotów punktualnych, opóźnionych, odwołanych oraz przekierowanych są publikowane w comiesięcznym raporcie „Air Travel Consumer Report” oraz w tym zestawie danych dotyczącym opóźnień i odwołań lotów z 2015 roku.

Pojedynczy wiersz zbioru danych reprezentuje pojedynczy lot. Zbiór danych zawiera informacje o czasie odlotu i przylotu, czasie opóźnienia, przyczynie opóźnienia, a także inne szczegóły dotyczące lotu.

## 2. Określenie celu eksploracji i kryteriów sukcesu

Celem eksploracji jest predykcja, czy lot będzie opóźniony, czy nie. Uznajemy, że lot jest opóźniony, jeśli atrybut `ARRIVAL_DELAY` jest większy niż 15. Dodatkowo chcemy zrozumieć, które czynniki mają największy wpływ na opóźnienia lotów.

Kryterium sukcesu jest osiągnięcie dokładności klasyfikacji opóźnienia na poziomie 85%. W naszym przypadku chcemy skupić się na maksymalizacji czułości modelu ponad swoistość, ponieważ błędne wykrycie opóźnienia nie jest tak istotne, jak przeoczenie rzeczywistego opóźnienia.

## 3. Charakterystyka zbioru

- Pochodzenie: [Kaggle](#)
- Liczba przykładów: 5819080
- Format: CSV (3 pliki: `flights.csv` - właściwy zbiór, `airports.csv` - informacje o lotniskach, `airlines.csv` - informacje o liniach lotniczych)
- Ilość zbiorów danych: 1

## 4. Opis atrybutów

Nazwa	Typ	Opis
<code>YEAR</code>	Numeryczny	Rok lotu
<code>MONTH</code>	Numeryczny	Miesiąc lotu
<code>DAY</code>	Numeryczny	Dzień miesiąca lotu
<code>DAY_OF_WEEK</code>	Numeryczny	Dzień tygodnia lotu
<code>AIRLINE</code>	Nominalny	Identyfikator linii lotniczej
<code>FLIGHT_NUMBER</code>	Numeryczny	Numer lotu
<code>TAIL_NUMBER</code>	Nominalny	Numer rejestracyjny samolotu
<code>ORIGIN_AIRPORT</code>	Nominalny	Kod IATA lotniska wylotu
<code>DESTINATION_AIRPORT</code>	Nominalny	Kod IATA lotniska przylotu
<code>SCHEDULED_DEPARTURE</code>	Numeryczny	Planowany czas odlotu (w formacie HHMM)
<code>DEPARTURE_TIME</code>	Numeryczny	Czas odlotu (w formacie HHMM)

Nazwa	Typ	Opis
DEPARTURE_DELAY	Numeryczny	Całkowite opóźnienie odlotu (w minutach)
TAXI_OUT	Numeryczny	Ilość minut spędzonych na kołowaniu przed odlotem
WHEELS_OFF	Numeryczny	Czas startu (w formacie HHMM)
SCHEDULED_TIME	Numeryczny	Planowany czas lotu (w minutach)
ELAPSED_TIME	Numeryczny	Całkowity czas lotu (w minutach) = $AIR\_TIME + TAXI\_IN + TAXI\_OUT$
AIR_TIME	Numeryczny	Czas lotu (w minutach)
DISTANCE	Numeryczny	Dystans lotu (w milach)
WHEELS_ON	Numeryczny	Czas lądowania (w formacie HHMM)
TAXI_IN	Numeryczny	Ilość minut spędzonych na kołowaniu po lądowaniu
SCHEDULED_ARRIVAL	Numeryczny	Planowany czas przylotu (w formacie HHMM)
ARRIVAL_TIME	Numeryczny	Czas przylotu (w formacie HHMM) = $WHEELS\_ON + TAXI\_IN$
ARRIVAL_DELAY	Numeryczny	Całkowite opóźnienie przylotu (w minutach) = $ARRIVAL\_TIME - SCHEDULED\_ARRIVAL$
DIVERTED	Prawda/Fałsz	Czy lot był przekierowany? (1-tak, 0-nie)
CANCELLED	Prawda/Fałsz	Czy lot był odwołany? (1-tak, 0-nie)
CANCELLATION_REASON	Nominalny	Przyczyna odwołania lotu (A - przewoźnik, B - pogoda, C - National Air System, D - Bezpieczeństwo), tylko dla odwołanych lotów
AIR_SYSTEM_DELAY	Numeryczny	Opóźnienie spowodowane przez system lotniczy (w minutach)
SECURITY_DELAY	Numeryczny	Opóźnienie spowodowane przez kontrole bezpieczeństwa (w minutach)
AIRLINE_DELAY	Numeryczny	Opóźnienie spowodowane przez przewoźnika (w minutach)
LATE_AIRCRAFT_DELAY	Numeryczny	Opóźnienie spowodowane przez samolot (w minutach)
WEATHER_DELAY	Numeryczny	Opóźnienie spowodowane przez pogodę (w minutach)

#### 4.1. Brakujące dane

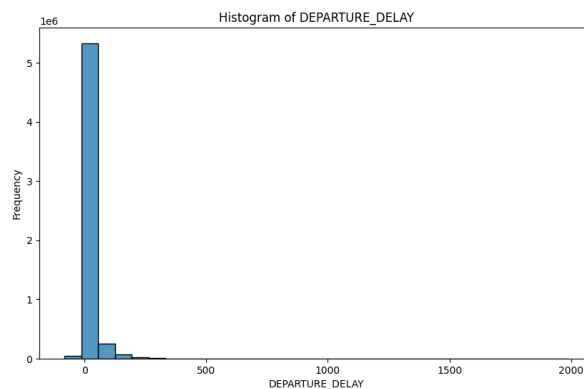
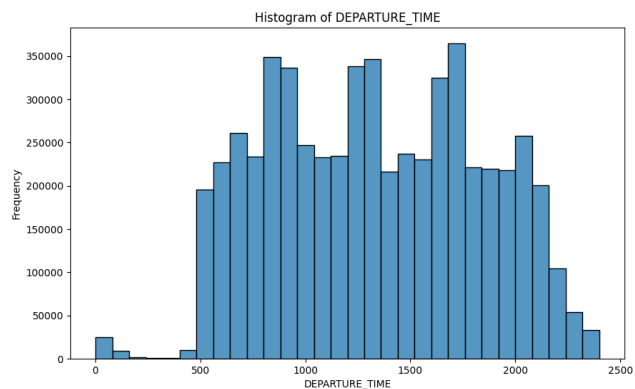
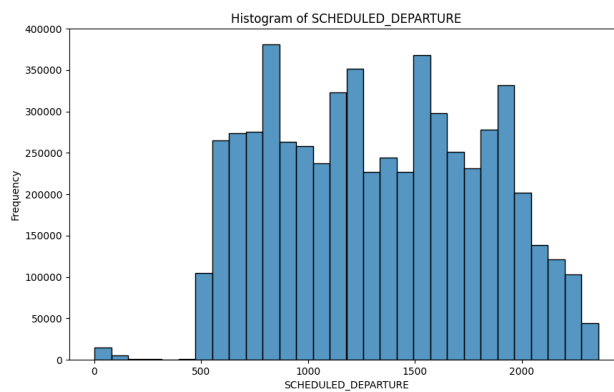
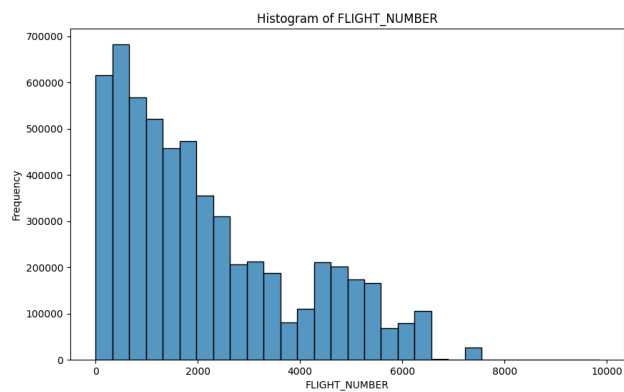
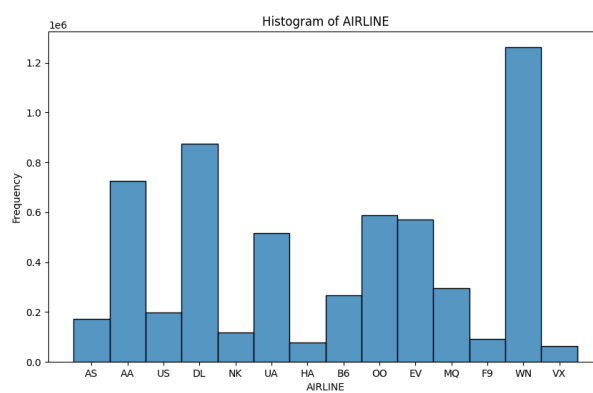
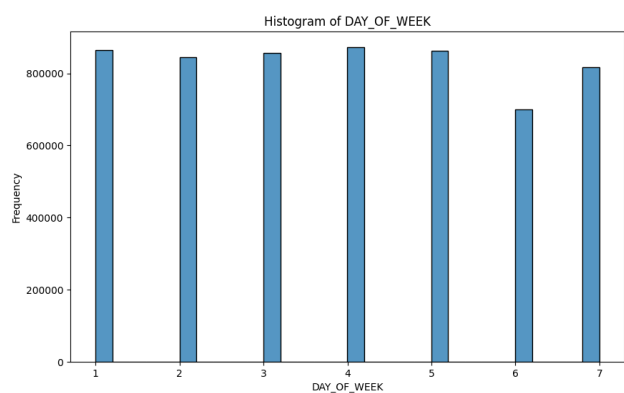
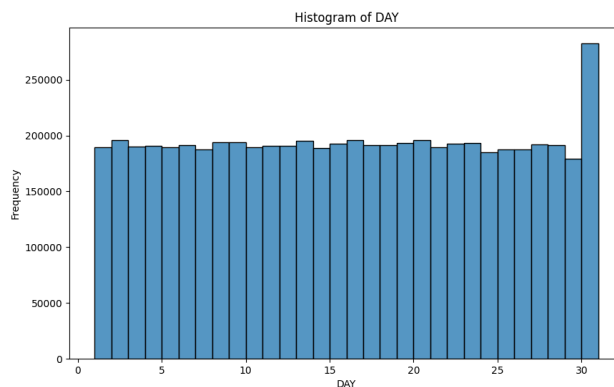
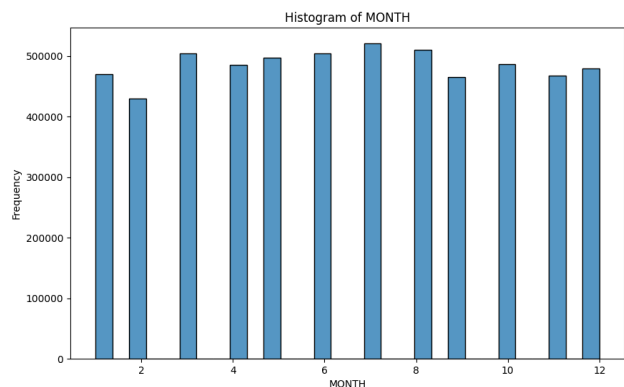
W zbiorze danych występowały nieznaczące braki, lecz stanowiły one nieznaczną część całego zbioru, który jest bardzo duży (ponad 5 milionów przykładów). Z uwagi na ten fakt, zostały one odfiltrowane podczas przetwarzania.

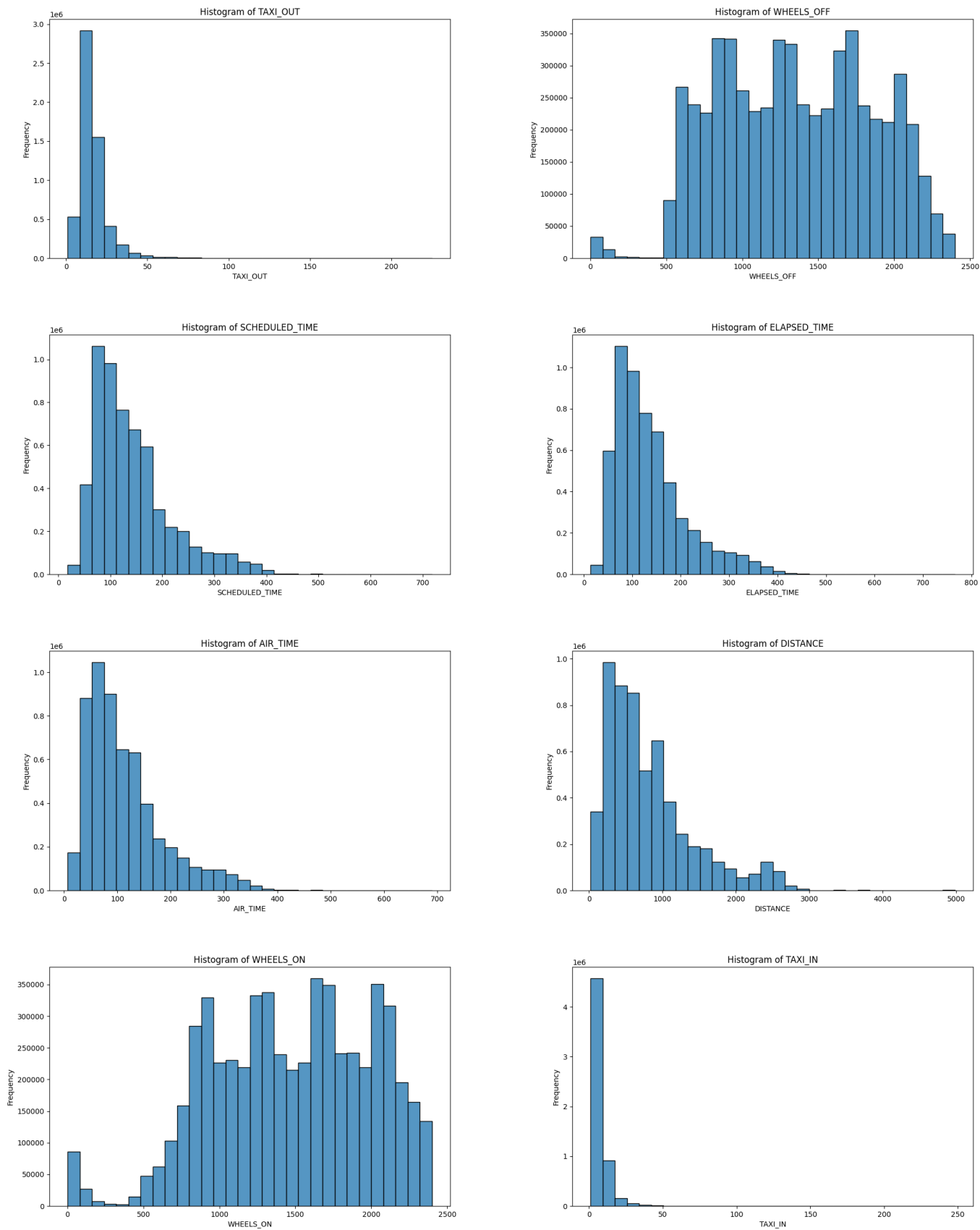
#### 4.2. Dane niezrozumiałe

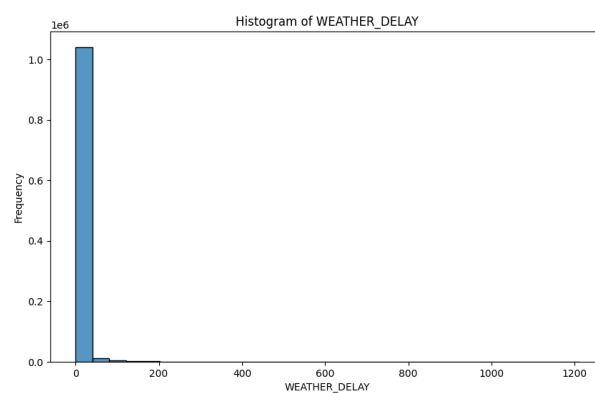
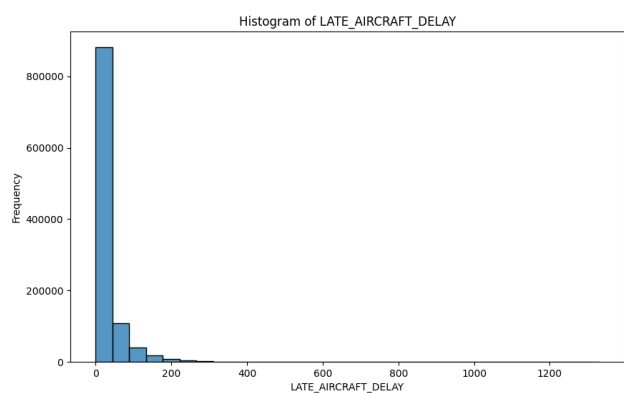
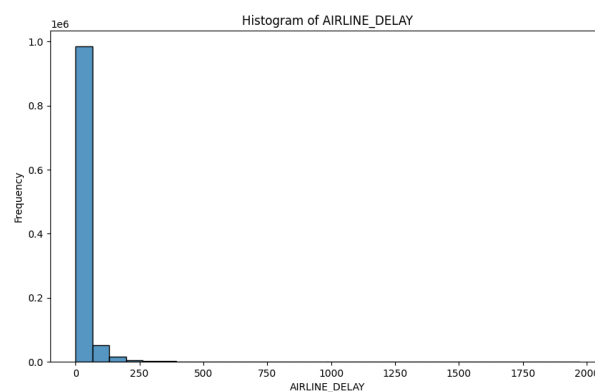
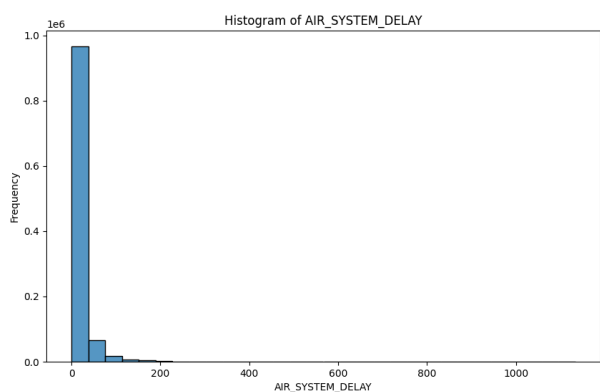
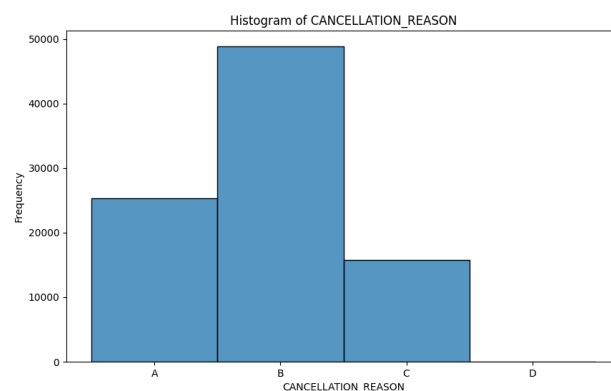
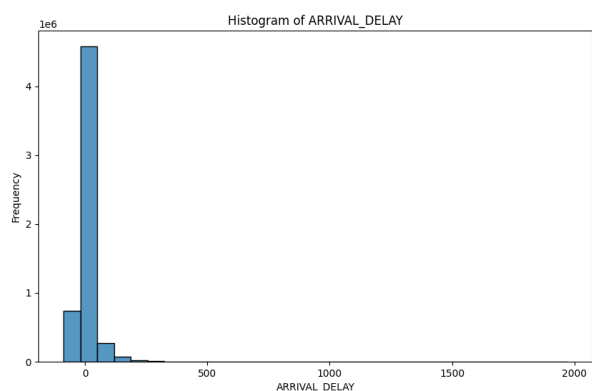
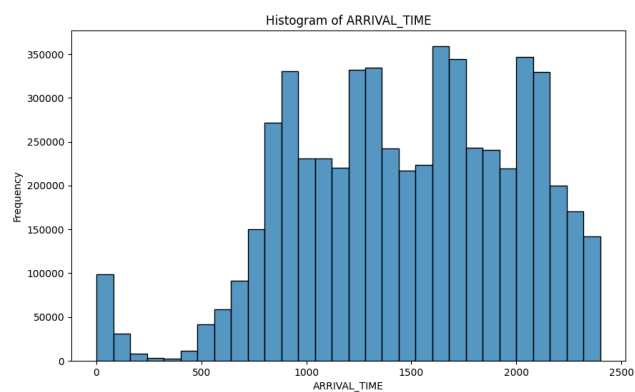
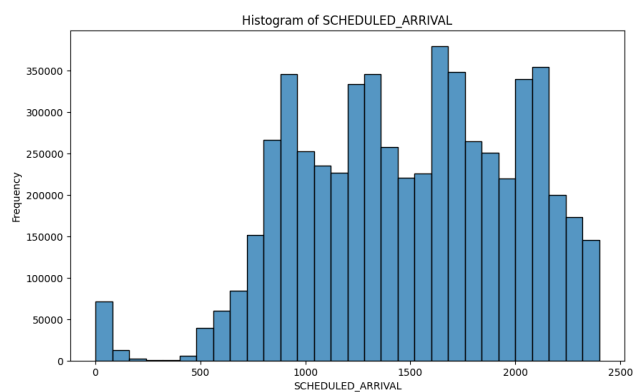
brak

## 5. Wynik eksploracyjnej analizy danych

### 5.1. Rozkłady wartości atrybutów



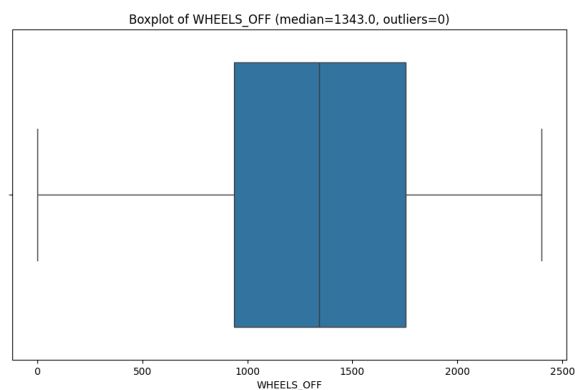
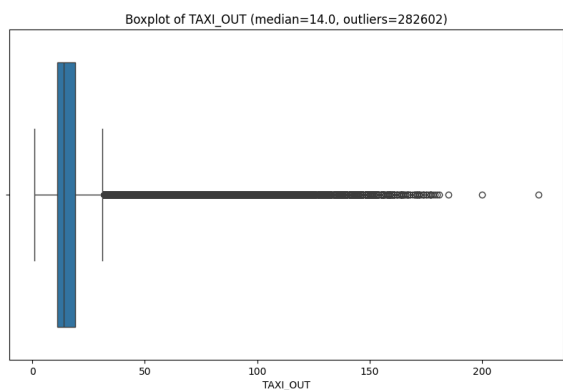
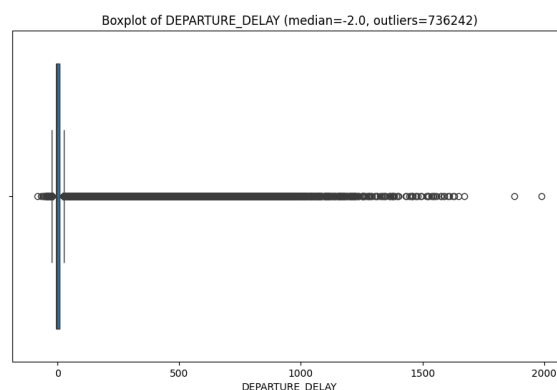
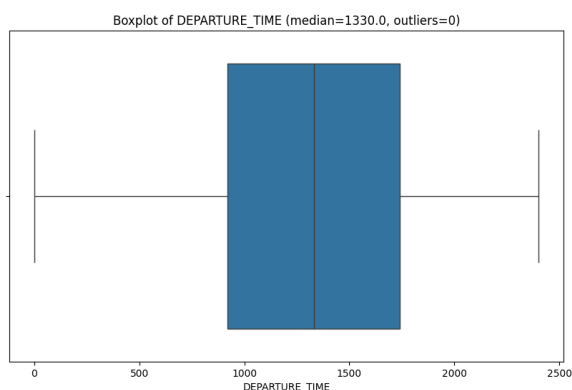
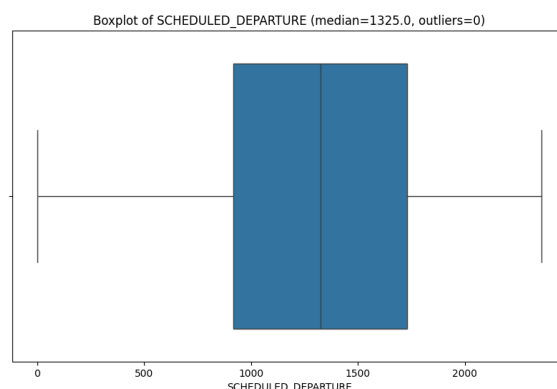
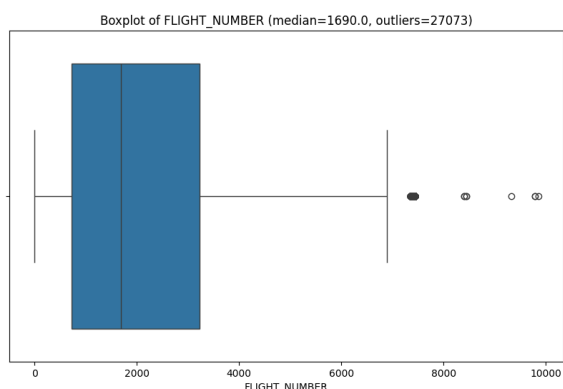


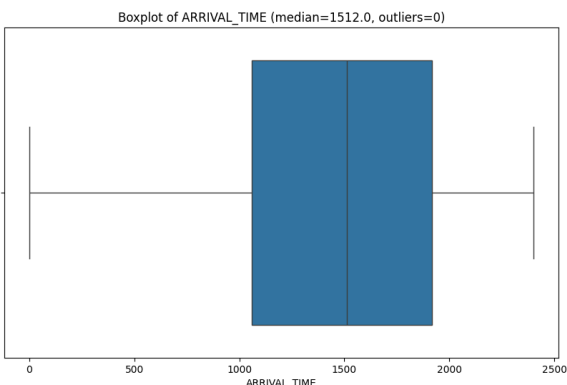
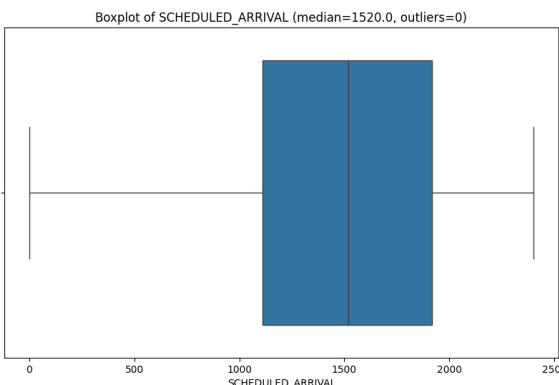
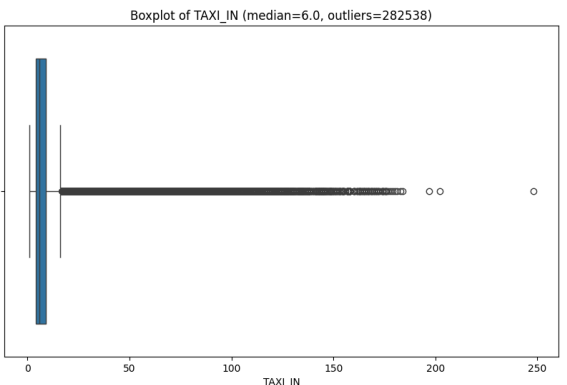
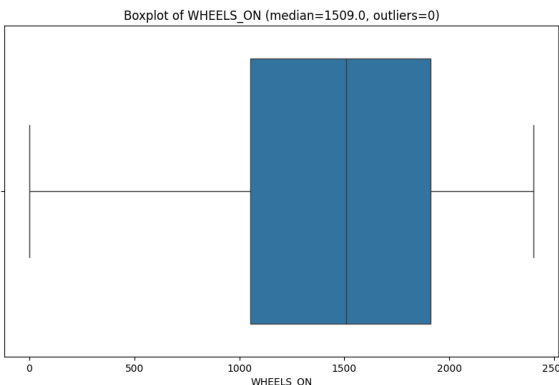
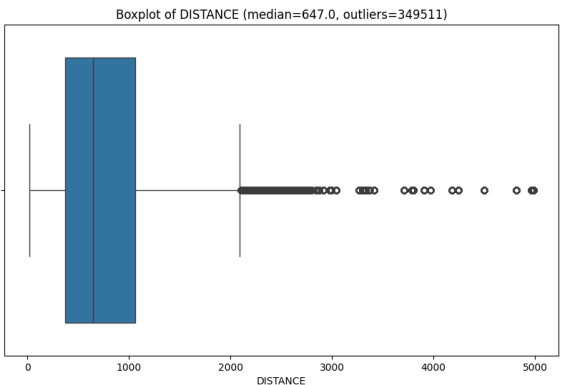
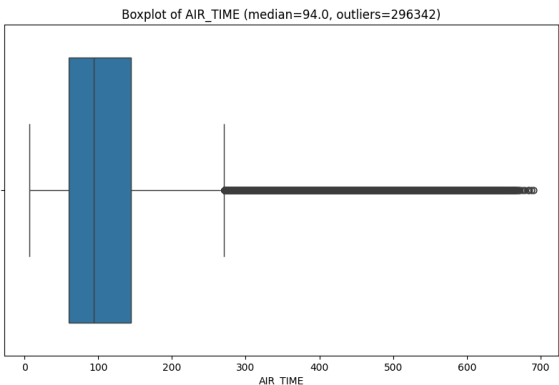
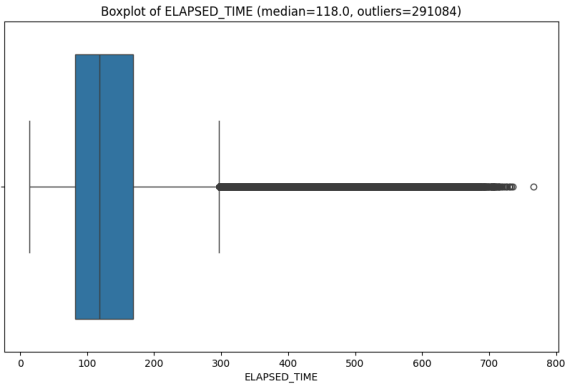
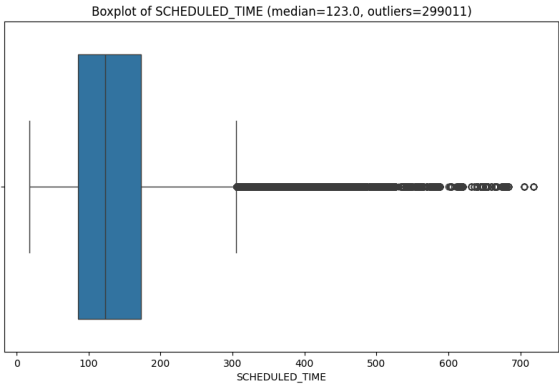


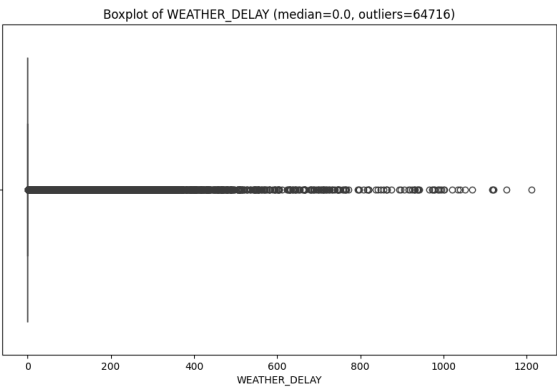
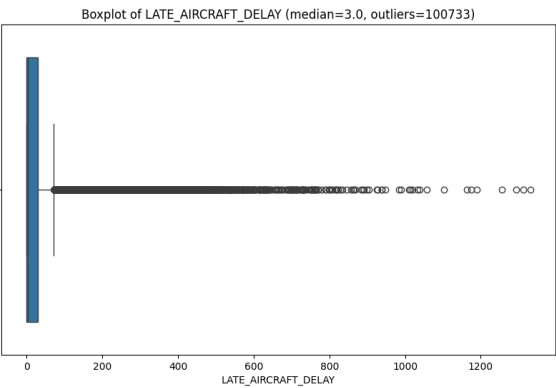
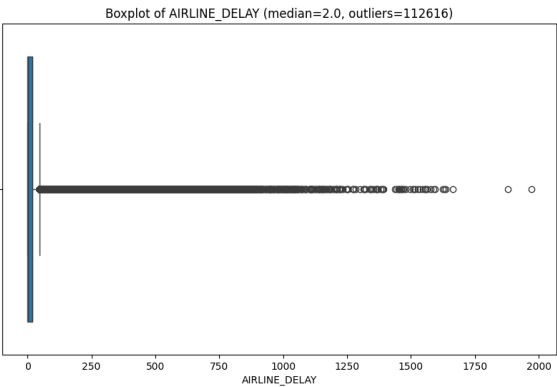
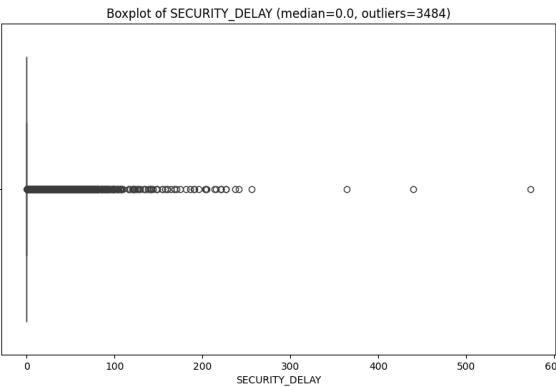
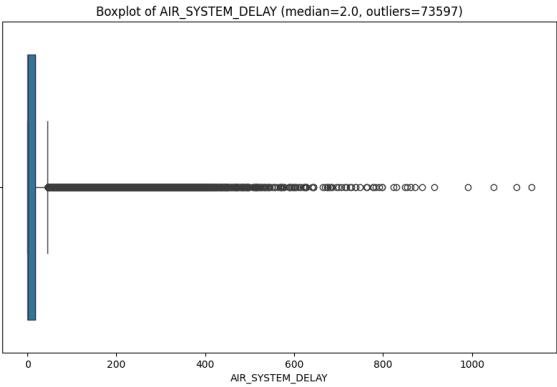
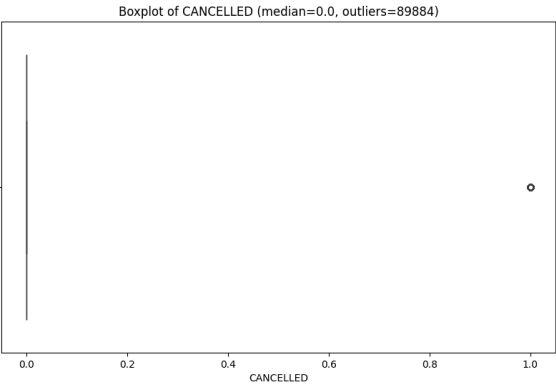
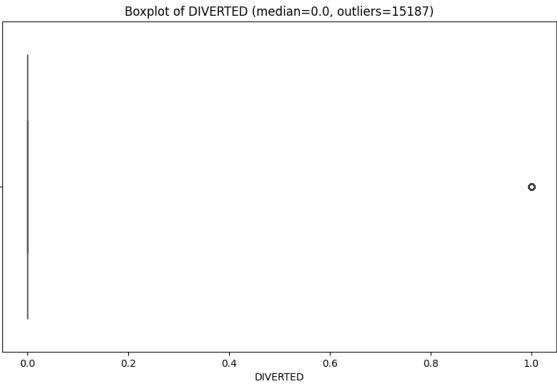
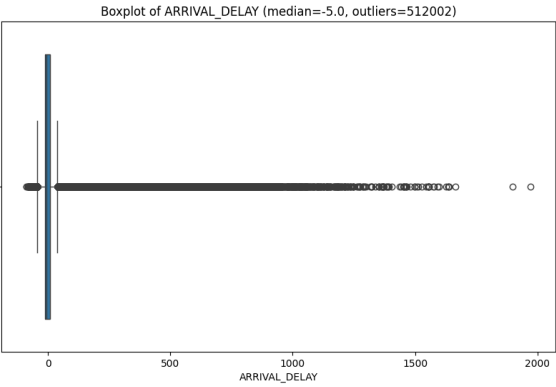
- W przypadku atrybutu `ARRIVAL_DELAY` zauważalna jest przewaga lotów punktualnych lub z niewielkim opóźnieniem (poniżej 15 minut) względem lotów znacząco opóźnionych. Klasy są niezbalansowane, co może wpłynąć na skuteczność modeli klasyfikacyjnych.

- Rozkłady większości atrybutów numerycznych, takich jak DEPARTURE\_DELAY, AIR\_TIME, TAXI\_OUT czy SCHEDULED\_TIME, nie przypominają rozkładu normalnego. Najczęściej obserwujemy rozkład prawoskośny – większość wartości skupia się w niższych przedziałach, a ogon rozkładu jest wydłużony w stronę wyższych wartości.
- Wysokie, rzadko występujące wartości w atrybutach takich jak WEATHER\_DELAY, AIRLINE\_DELAY oraz LATE\_AIRCRAFT\_DELAY mogą wskazywać na zdarzenia nietypowe, takie jak intensywne burze, problemy techniczne lub opóźnienia łańcuchowe wynikające z wcześniejszych lotów. Tego typu przypadki mają istotne znaczenie dla analizy przyczyn opóźnień i mogą być kluczowe przy budowie predykcyjnych modeli.
- Wartości atrybutu DISTANCE rozkładają się nierównomiernie – większość lotów odbywa się na krótkich i średnich dystansach, co znajduje odzwierciedlenie w rozkładzie. Długodystansowe loty są mniej liczne.

## 5.2. Punkty oddalone







Nazwa	Mediana	Punkty oddalone
AIR_SYSTEM_DELAY	2	73 597
AIR_TIME	94	296 342



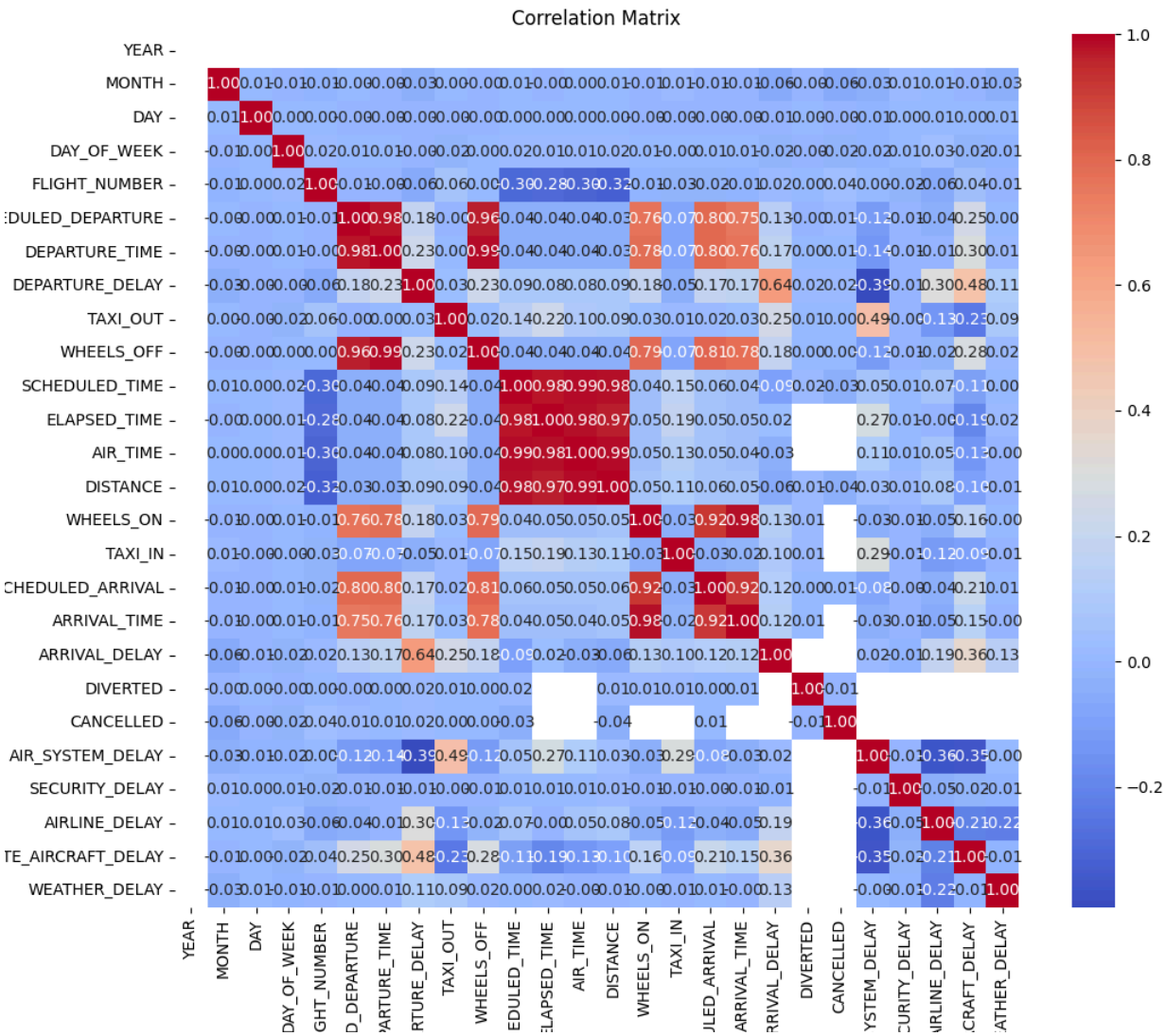
Nazwa	Mediana	Punkty oddalone
AIRLINE_DELAY	2	112 616
ARRIVAL_DELAY	-5	512 002
ARRIVAL_TIME	1 512	0
CANCELLED	0	89 884
DAY_OF_WEEK	4	0
DAY	16	0
DEPARTURE_DELAY	-2	736 242
DEPARTURE_TIME	1 330	0
DISTANCE	647	349 511
DIVERTED	0	15 187
ELAPSED_TIME	118	291 084
FLIGHT_NUMBER	1 690	27 073
LATE_AIRCRAFT_DELAY	3	100 733
MONTH	7	0
SCHEDULED_ARRIVAL	1 520	0
SCHEDULED_DEPARTURE	1 325	0
SCHEDULED_TIME	123	299 011
SECURITY_DELAY	0	3 484
TAXI_IN	6	282 538
TAXI_OUT	14	282 602
WEATHER_DELAY	0	64 716
WHEELS_OFF	1 343	0
WHEELS_ON	1 509	0
YEAR	2015	0

Na podstawie przedstawionych median oraz liczby punktów oddalonych dla poszczególnych atrybutów można sformułować następujące wnioski:

- Dla większości atrybutów numerycznych, takich jak AIR\_SYSTEM\_DELAY, AIRLINE\_DELAY, DEPARTURE\_DELAY czy LATE\_AIRCRAFT\_DELAY, mediana wynosi 0, 2 lub wartości bliskie zeru. Oznacza to, że typowy lot nie doświadcza istotnych opóźnień z tych przyczyn, a większość lotów przebiega zgodnie z planem.
- Wysoka liczba punktów oddalonych (np. 736 242 dla DEPARTURE\_DELAY, 299 011 dla SCHEDULED\_TIME, 296 342 dla AIR\_TIME) wskazuje na obecność nietypowych, ekstremalnych przypadków w danych. Takie wartości mogą być wynikiem wyjątkowych zdarzeń, np. bardzo dużych opóźnień, awarii lub specyficznych tras.
- Mediana opóźnienia przylotu (ARRIVAL\_DELAY) jest ujemna (-5), co sugeruje, że ponad połowa lotów przylatuje przed planowanym czasem lub z minimalnym opóźnieniem.
- Atrybuty binarne, takie jak CANCELLED czy DIVERTED, mają medianę 0, co oznacza, że większość lotów nie jest odwoływana ani przekierowywana.

- Dla atrybutów czasowych (np. ARRIVAL\_TIME, DEPARTURE\_TIME, WHEELS\_ON, WHEELS\_OFF) mediany odpowiadają typowym godzinom operacji lotniczych, a brak punktów oddalonych sugeruje, że wartości te są stabilne.
- Wysoka liczba punktów oddalonych w atrybutach związanych z czasem trwania lotu (AIR\_TIME, ELAPSED\_TIME, SCHEDULED\_TIME) oraz dystansem (DISTANCE) odzwierciedla zróżnicowanie tras – od krótkich po bardzo długie loty.

### 5.3. Macierz korelacji



Istnieje silna korelacja pomiędzy atrybutami WHEELS\_OFF i ARRIVAL\_TIME (0.78), co pokazuje że większość opóźnień wynika z opóźnieniami na lądzie. Same przeloty są punktualne. Podobne wnioski można wysnuć w przypadku atrybutów SECURITY\_DELAY i WEATHER\_DELAY które prawie nie wykazują korelacji z całkowitym opóźnieniem co może sugerować że czasy odprawy i pogoda są niewielkim czynnikiem wpływającym na opóźnienia.

## 6. Podsumowanie

Na początku przeprowadzono analizę rozkładów atrybutów, z której wynika, że większość zmiennych numerycznych nie ma rozkładu normalnego — najczęściej przyjmują one postać rozkładów

prawoskośnych. Dodatkowo klasy zmiennej celu (czy lot jest opóźniony) są niezbalansowane — znacznie więcej lotów kończy się punktualnie lub z niewielkim opóźnieniem. W związku z nienormalnością rozkładów zastosowano współczynnik korelacji rang Spearmana do analizy zależności między zmiennymi.

Analiza korelacji wykazała istnienie silnych związków pomiędzy niektórymi grupami atrybutów. Szczególnie silna korelacja występuje między momentem startu (WHEELS\_OFF) a czasem przylotu (ARRIVAL\_TIME), co sugeruje, że opóźnienia przylotu wynikają przede wszystkim z opóźnień przy starcie, a sam czas przelotu jest relatywnie stabilny. Z kolei atrybuty takie jak WEATHER\_DELAY czy SECURITY\_DELAY wykazują bardzo niską korelację z opóźnieniem przylotu, co może świadczyć o ich mniejszym znaczeniu w kontekście predykcji.

Oceniono również jakość danych. Występuje bardzo duża liczba punktów odstających w wielu atrybutach (np. DEPARTURE\_DELAY, AIR\_TIME, SCHEDULED\_TIME), co może wskazywać na nietypowe sytuacje operacyjne, takie jak intensywne warunki pogodowe, awarie techniczne czy długodystansowe trasy. Mimo to, z racji potencjalnej informacyjności takich obserwacji, zdecydowano się ich nie usuwać, ponieważ mogą mieć istotne znaczenie dla budowy modelu predykcyjnego.

Na podstawie powyższej analizy można stwierdzić, że dane są wystarczająco dobrej jakości, by możliwe było osiągnięcie postawionego celu eksploracji – predykcji opóźnień przylotów z dokładnością co najmniej 85%. Szczególna uwaga zostanie położona na maksymalizację czułości modelu, aby ograniczyć liczbę przypadków, w których rzeczywiste opóźnienie nie zostanie wykryte.