

# Projekt IUM - 2021Z - Zad.3

Aleksander Garlikowski  
Cezary Moczulski

## **Problem biznesowy:**

Przewidywanie czasu transportu zakupionego produktu.

## **Zadanie modelowania:**

Regresja czasu dostawy na podstawie dnia tygodnia zamówienia, miasta docelowego oraz firmy transportowej obsługującej zamówienie.

## **Biznesowe kryterium sukcesu:**

Przekazywanie klientom przewidywań na temat zakresu czasu dostawy, dokładniejszych i bardziej wiarygodnych niż rozsądne przewidywanie, że produkt zostanie dostarczony w zakresie od 24 do 72 godzin od momentu zakupu (na podstawie dostępnych danych sprawdza się ono w ~79% przypadków).

## **Analityczne kryterium sukcesu:**

Uzyskanie przewidywanego zakresu czasu węższego niż 48 godzin, przy celności wyższej niż 90% (żeby można było uznać przewidywanie za znacznie bardziej wiarygodne i zostawić margines błędu dla wdrożenia modelu), na posiadanych danych.

90% celności powinno być wynikiem osiągalnym, ponieważ nawet przewidywanie polegające na podawaniu średniej dla każdej z kombinacji miasta, firmy kurierskiej i dnia tygodnia osiągać ponad 90% skuteczności z oknem czasowym o szerokości 48 godzin.

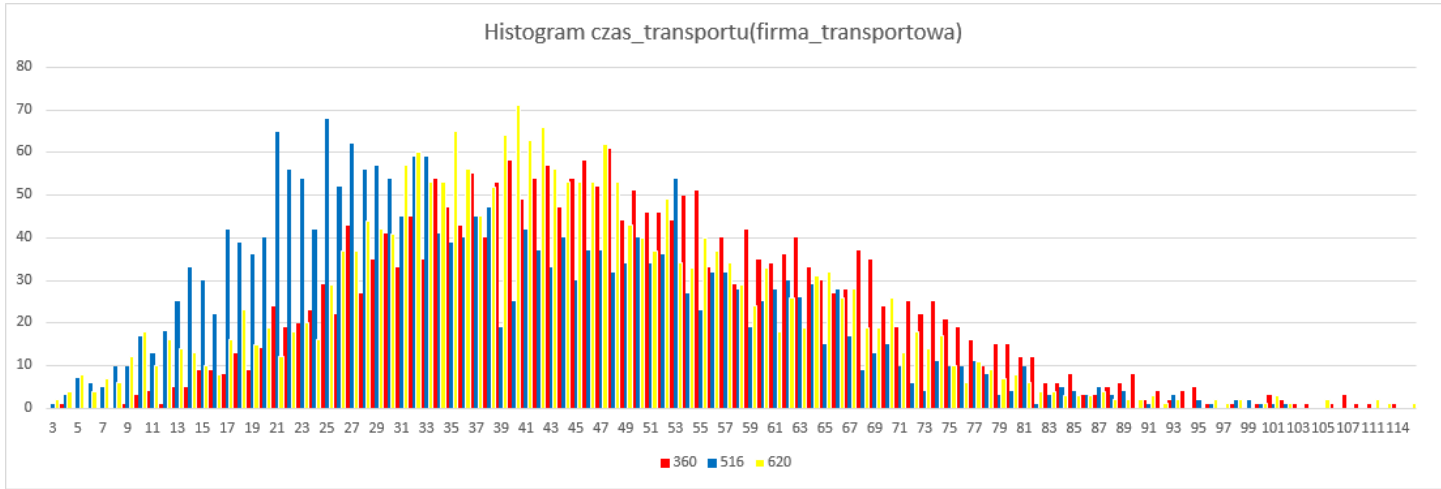
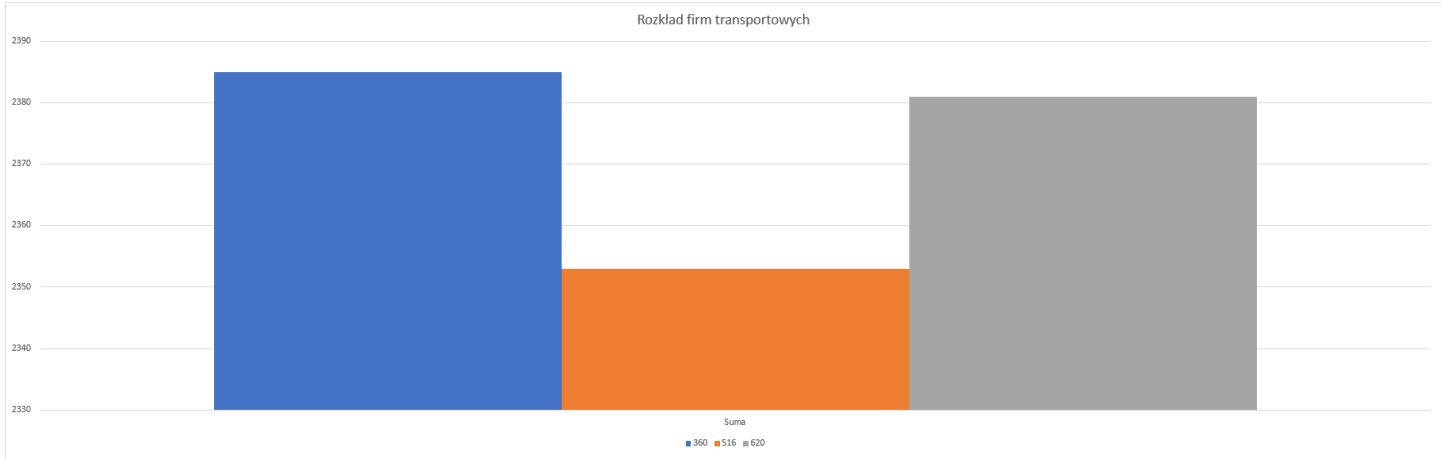
# Analiza danych:

Wybrane dane:

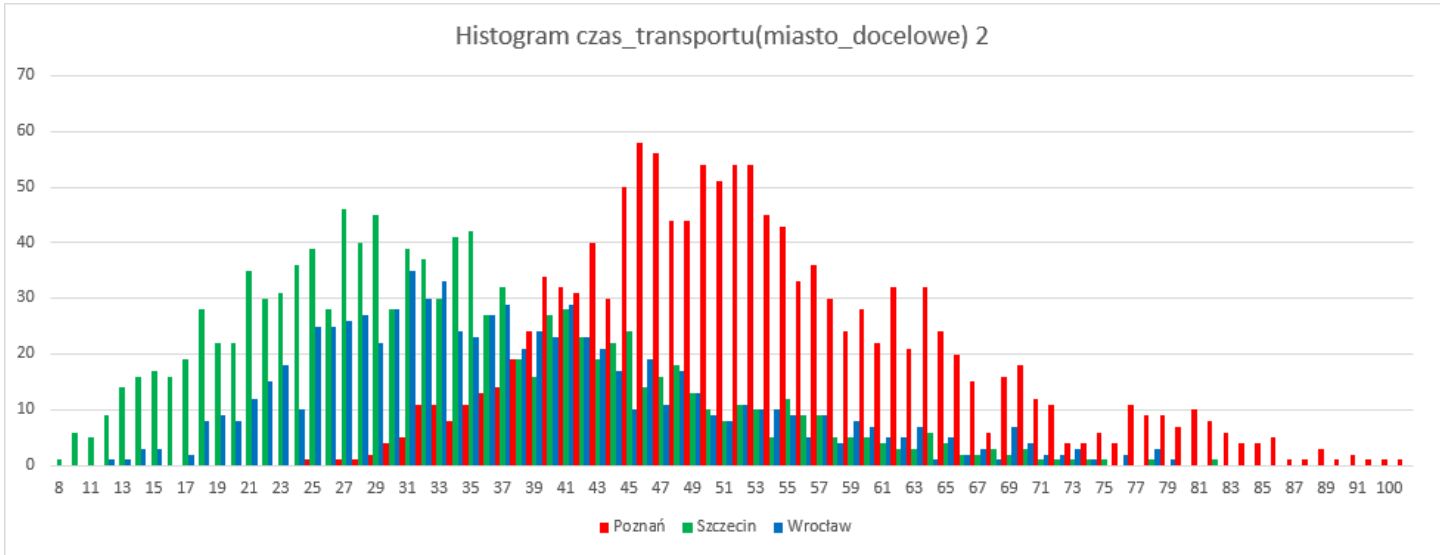
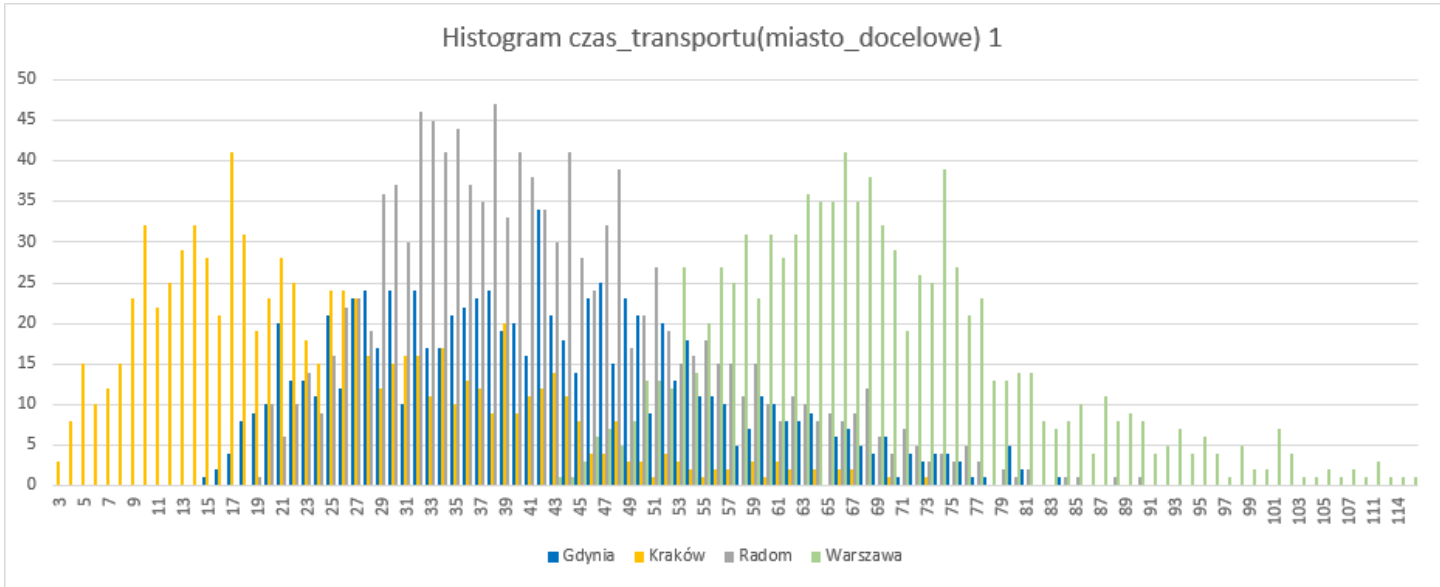
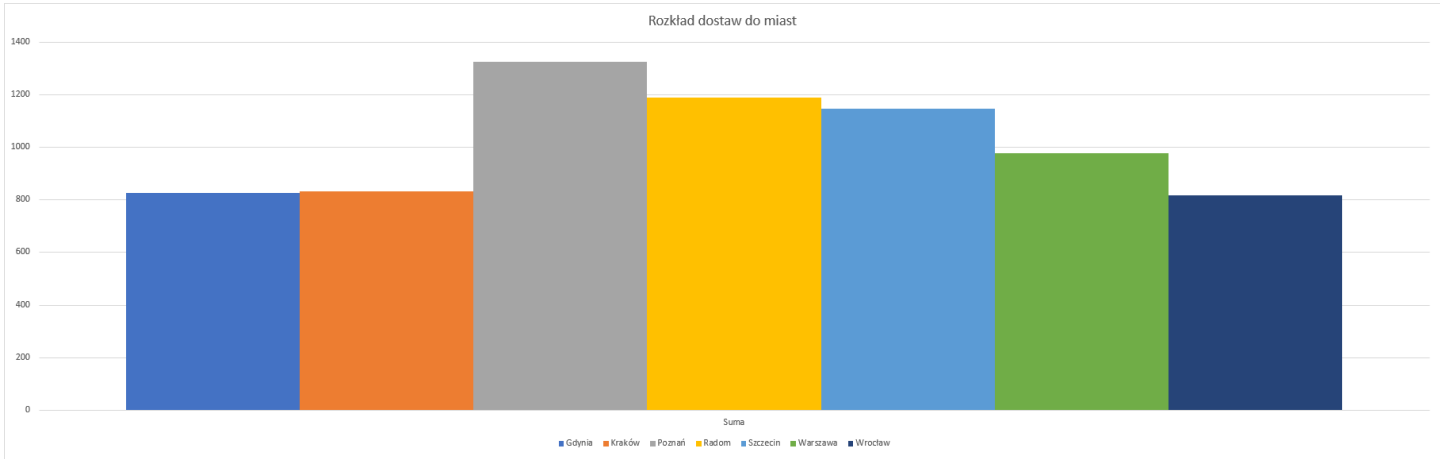
- Firma kurierska
- Miasto docelowe
- Dzień tygodnia złożenia zamówienia

Poniżej znajdują się wykresy rozkładu wystąpień atrybutów w danych i rozkłady czasu transportu (zmiennej celu) od tych atrybutów.

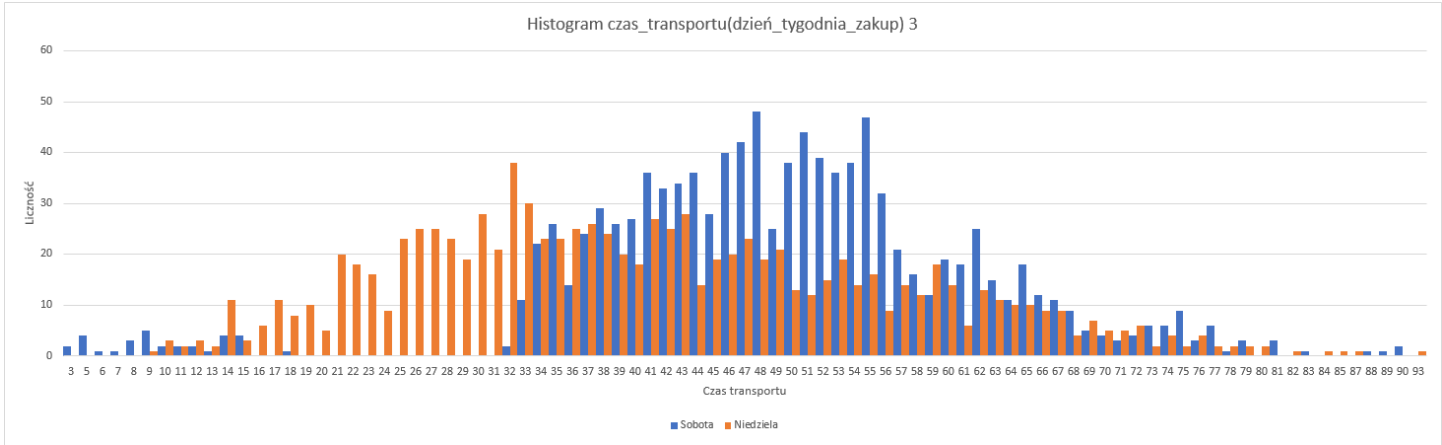
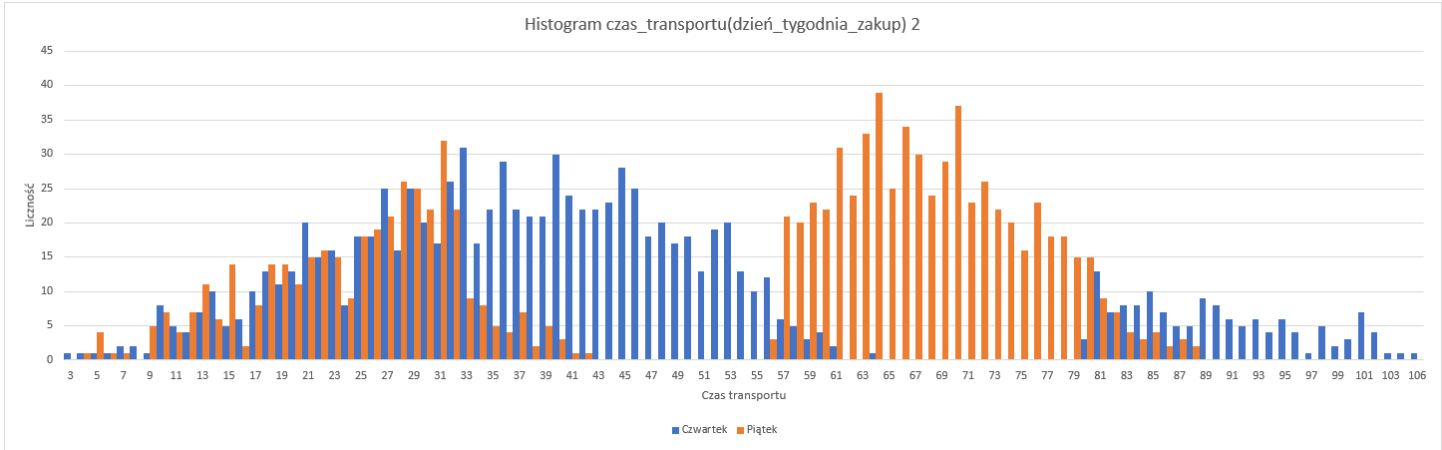
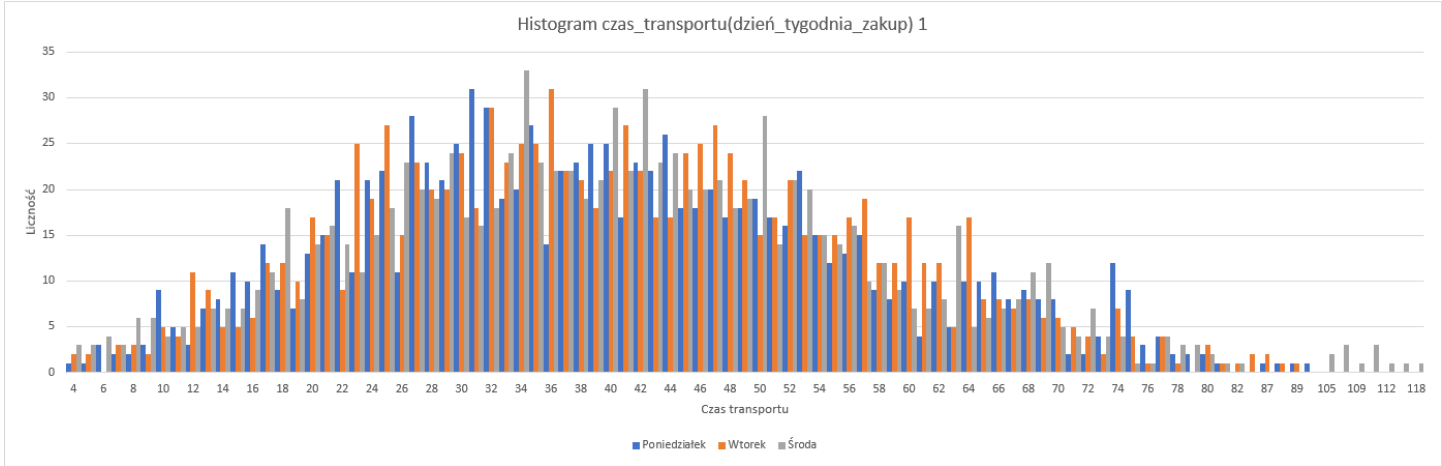
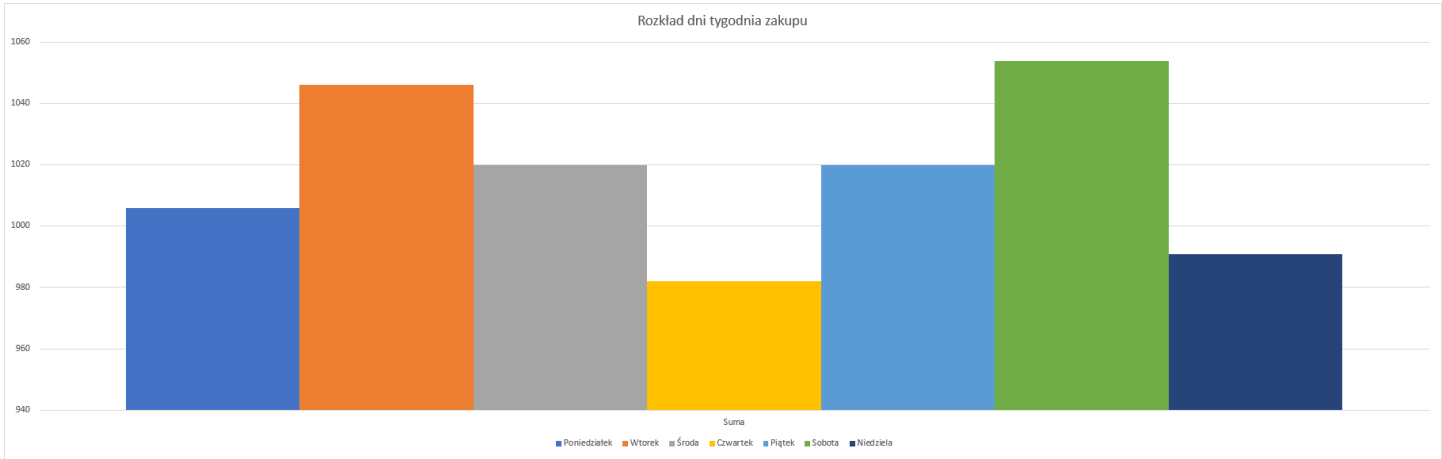
Firma kurierska:



Miasto:



Dzień tygodnia:



Odrzucone dane:

Odrzucone na podstawie wiedzy dziedzinowej:

- Imię i nazwisko klienta
- Nazwa produktu
- Cena produktu
- Przecena

Odrzucone, ponieważ nie niosą znaczącej informacji o zmiennej celu:

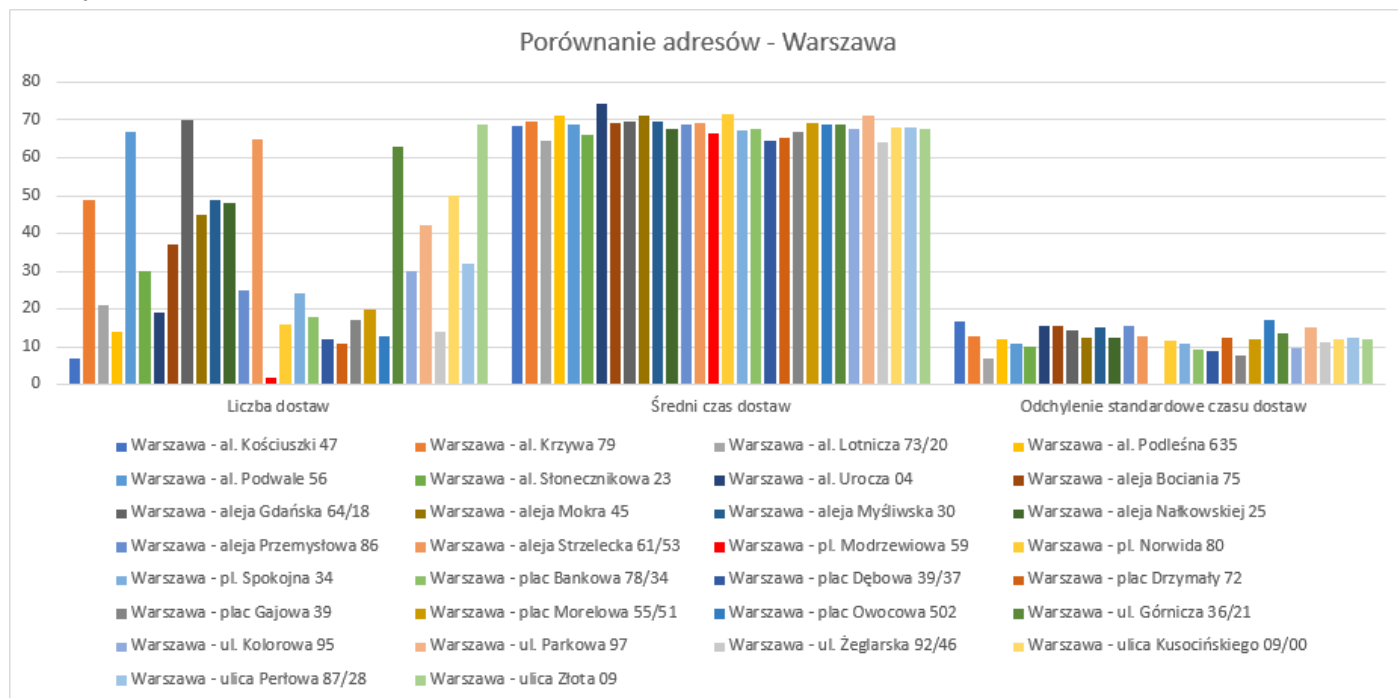
- Nazwy ulic (adresy)
- Kategoria produktu (gabaryty najwidoczniej nie mają znaczenia)
- Miesiąc
- Dokładna data (np. święta)
- Godzina

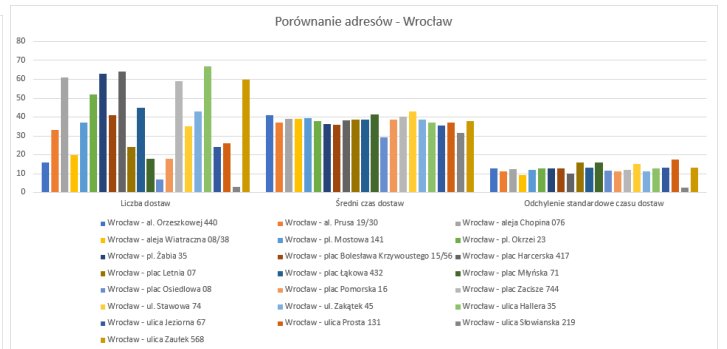
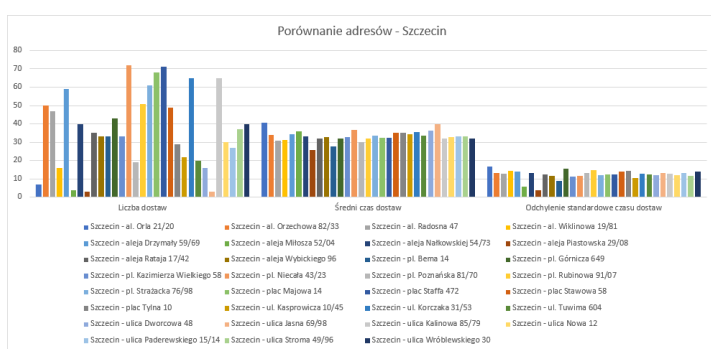
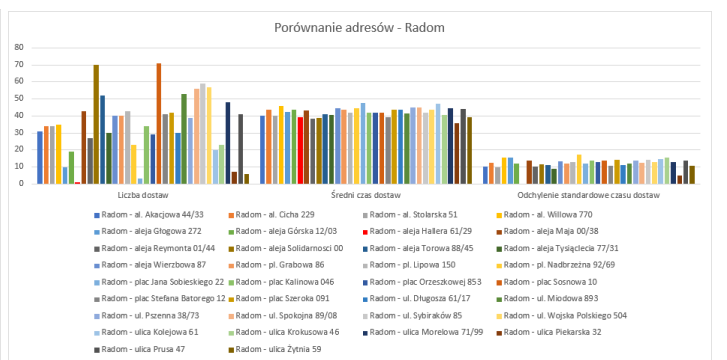
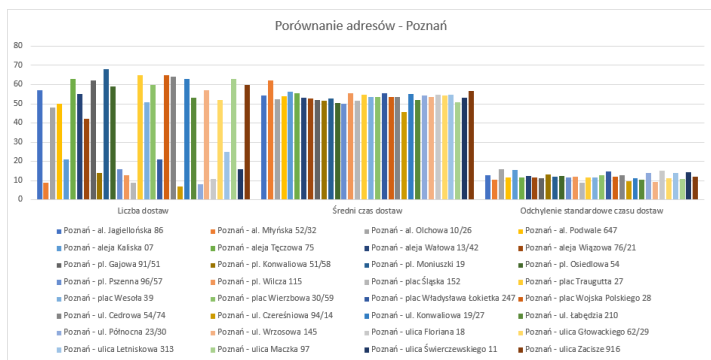
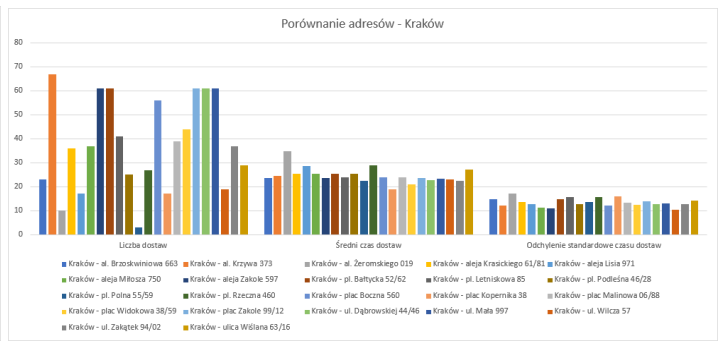
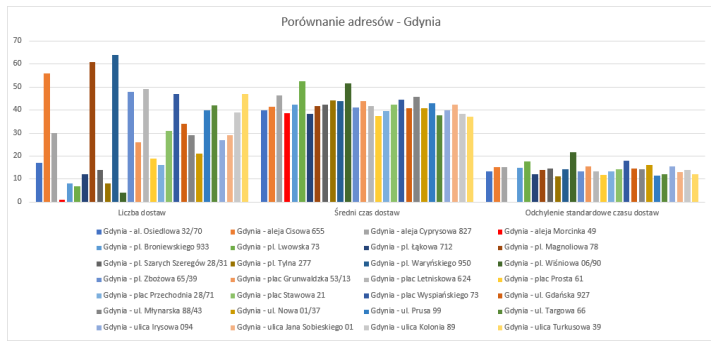
Ponieważ nie mamy danych pochodzących z innych lat, nie jesteśmy w stanie stwierdzić czy rok niesie ze sobą istotne informacje. Z daty uwzględniamy więc jedynie dzień tygodnia.

Po przejrzeniu danych stwierdziliśmy, że rozkład czasu transportu (zmiennej celu) dla konkretnej wartości każdego z odrzuconych atrybutów jest bardzo podobny do rozkładów dla jego pozostałych wartości.

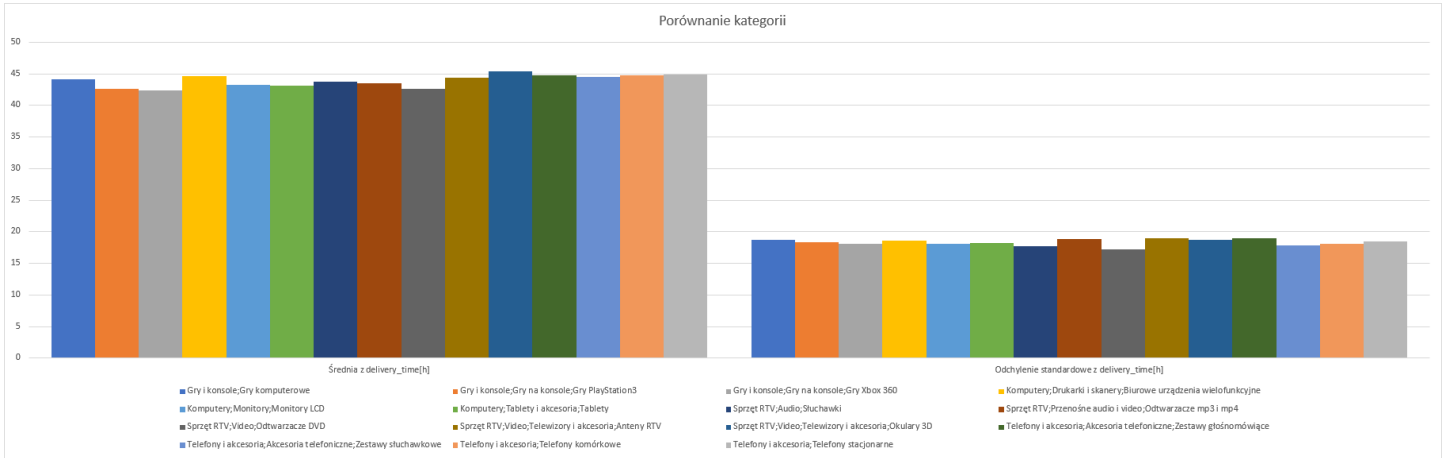
Poniżej są wykresy średnich i odchyłeń standardowych odrzuconych atrybutów.

Adresy w ramach miasta:

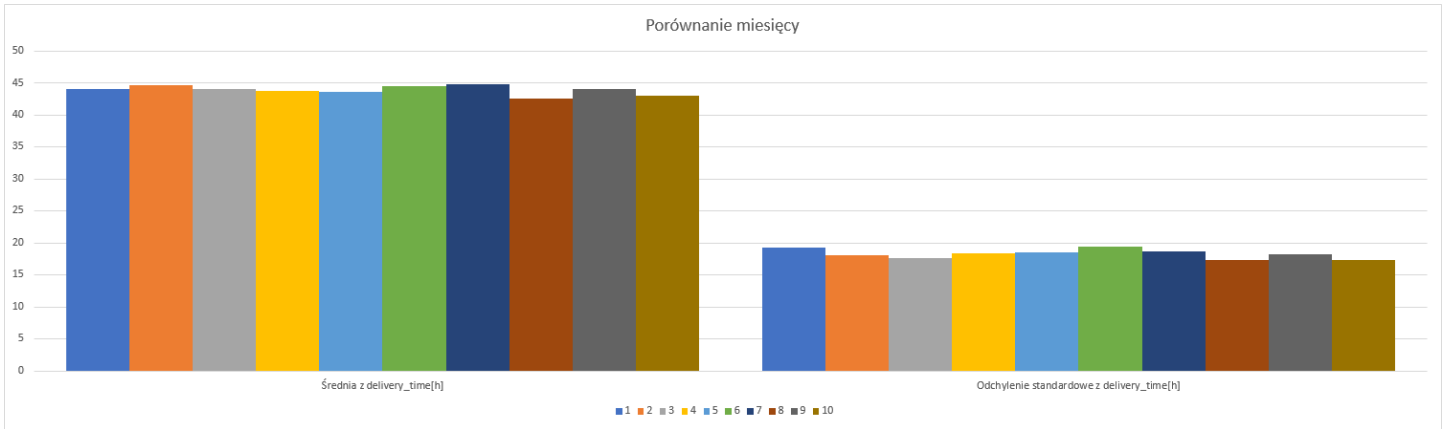




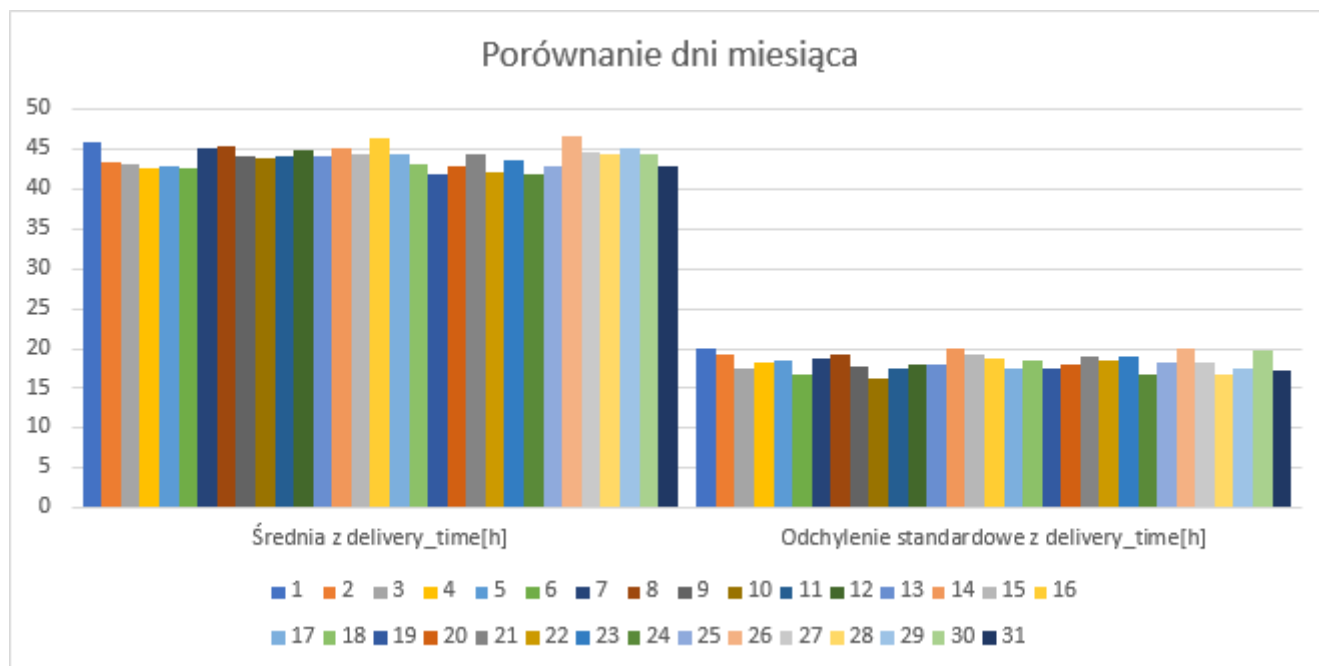
Kategoria produktu:



Miesiąc:



Dzień miesiąca:



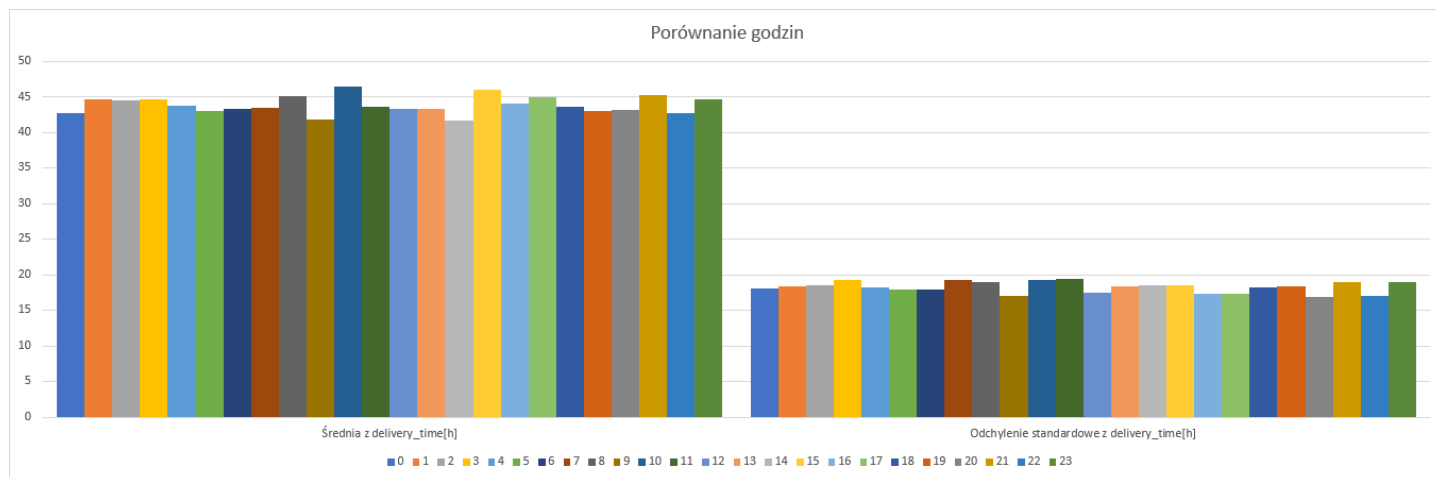
Tutaj lekka rozbieżność średnich wynika z wpływu dnia tygodnia na czas transportu.

Specjalne dni, np. święta:

W święta przesyłki też są dostarczane i średnia czasów dostaw dla świąt jest bardzo zbliżona do ogólnej.

Nie sprawdzaliśmy rozkładu dla każdego dnia osobno, ale nie natknęliśmy się na żadne odstające dni lub większe skupiska dłuższych dostaw w pewnych okresach czasu.

Godzina zakupu:





## **Założenia:**

- Zakupione produkty nie ulegają “zagubieniu” w trakcie transportu, zawsze docierają do klienta.
- Firmy kurierskie nie pracują w niedziele.
- Dostawy odbywają się święta tak samo jak w normalne dni.
- Dostawy odbywają się jedynie w godzinach 8.00-19.00, niezależnie od firmy transportowej.
- Model będzie przewidywał czasy transportu dla ograniczonej liczby miast i firm transportowych
- Dane zawsze będą kompletne, tak jak te, które otrzymaliśmy (niekompletne rekordy mogą być odrzucane)

## **Kompletność i reprezentatywność danych:**

W danych nie stwierdzono brakujących atrybutów ani wyraźnych sprzeczności (wszystkie rekordy zawierają wszystkie pola, czasy dostaw nie są ujemne).

Do pełnego obrazu brakuje jednak dostaw dokonanych od końca października do końca roku. Mimo, że nie uwzględniamy ich w modelu, mogłoby się okazać, że powinniśmy, ponieważ akurat w tym czasie coś się zmienia (np. pojawiają się opóźnienia związane z liczbą zamówień przed świętami).

Żeby model dostosowywał się lepiej do zjawisk zmieniających się na przestrzeni lat, przydatna byłaby również chociaż niewielka ilość, najlepiej równo rozłożonych, danych z kilku różnych lat ubiegłych.

Ostatecznie jednak dane, dla podanej sytuacji, uznajemy za wystarczające do skonstruowania modelu spełniającego wyznaczone przez nas kryteria. Model może jednak nie być w stanie dostosować się do zmian w ciągu kolejnych lat i będzie ograniczony do tych miast i firm kurierskich, które występują obecnie w danych.