

Raport: Budowanie i Weryfikacja Hipotez Dotyczących Predykcji Czasu Dostawy

1. Wstęp

Obecny algorytm przewiduje czasy dostawy na podstawie globalnej średniej ze wszystkich danych dla każdego zamówienia stosowana jest ta sama wartość. Takie podejście jest bardzo uproszczone i nie oddaje złożoności rzeczywistego procesu dostawy. W niniejszym raporcie przedstawiamy nasze założenia oraz metodologię, którą wykorzystaliśmy do analizy dostępnych danych. Wyciągamy wnioski, wskazując, jakie kroki można podjąć, aby prognozy były bardziej precyzyjne, a projektowany czas lepiej odpowiadał rzeczywistym warunkom.

2. Założenia i Metodologia

Założenia:

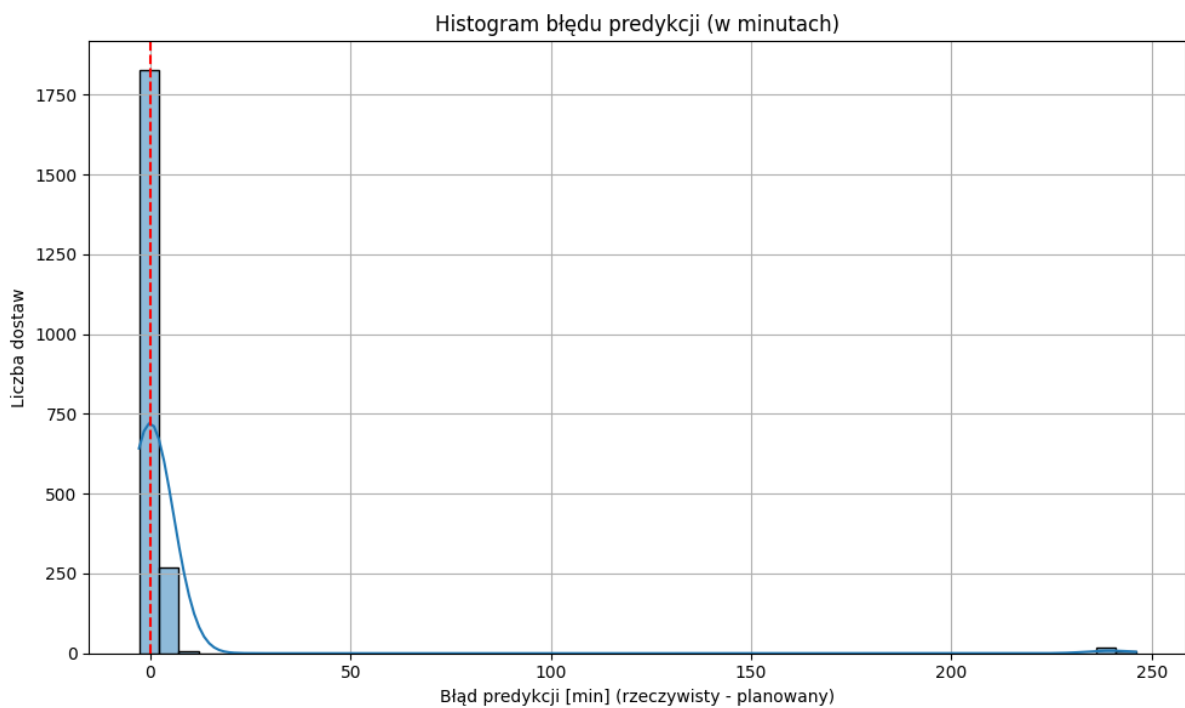
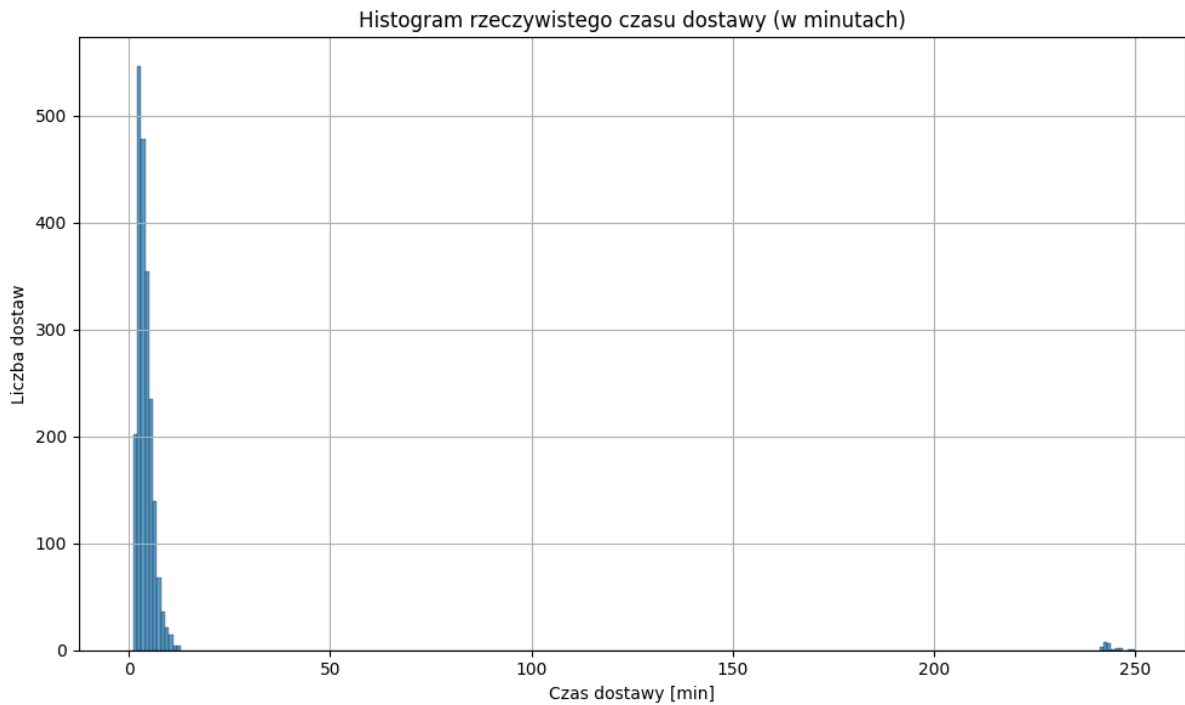
- **Brak założeń o jakości danych:** Nie przyjmujemy, że dane są idealne, przed rozpoczęciem analizy usuwaliśmy rekordy z brakującymi lub podejrzanymi (np. ujemnymi) wartościami.
- **Obecny algorytm jest bardzo uproszczony:** Planowany czas dostawy jest ustalany jako jedna, stała wartość uzyskaną jako średnia ze wszystkich dotychczasowych zamówień. To podejście nie różnicuje warunków dla zamówień wykonywanych w różnych sektorach, o różnych porach dnia czy o różnej wielkości zamówienia.
- **Brak danych o typie budynku:** Nie posiadamy informacji, czy dostawa odbywa się do domu jednorodzinnego, czy do budynku wielorodzinnego, dlatego szukamy proxy (np. sektor, liczba produktów, waga zamówienia), które mogą wskazywać na różnice w warunkach dostawy.

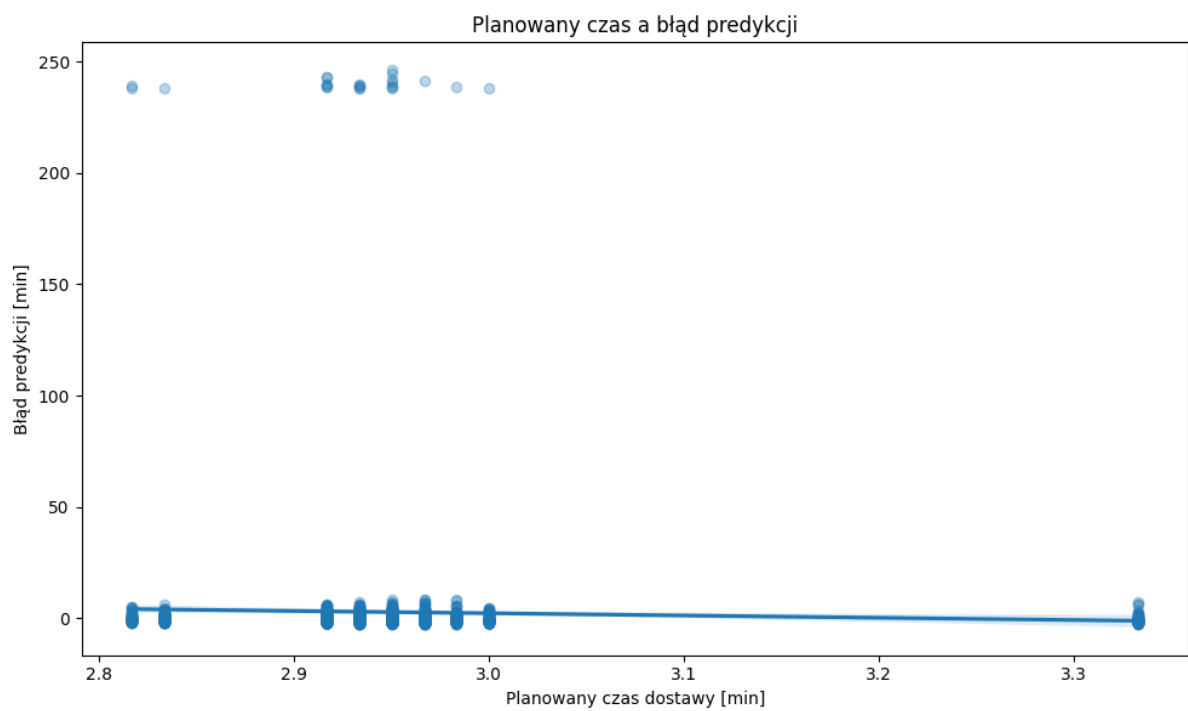
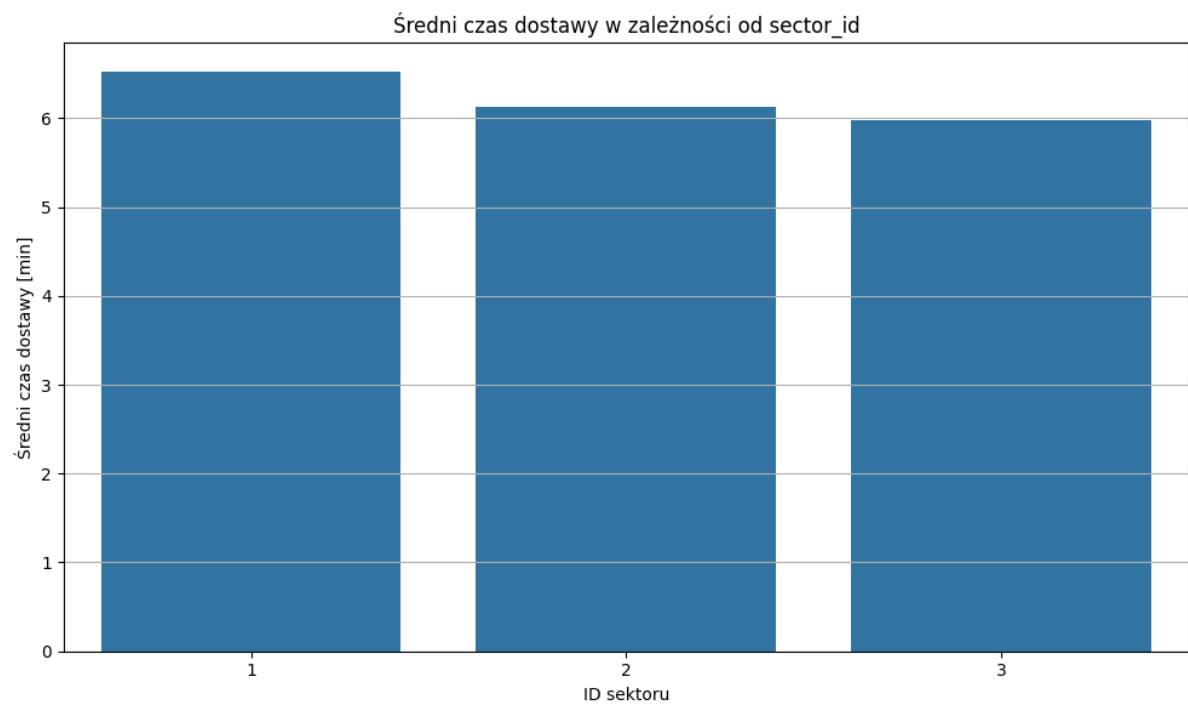
Metodologia:

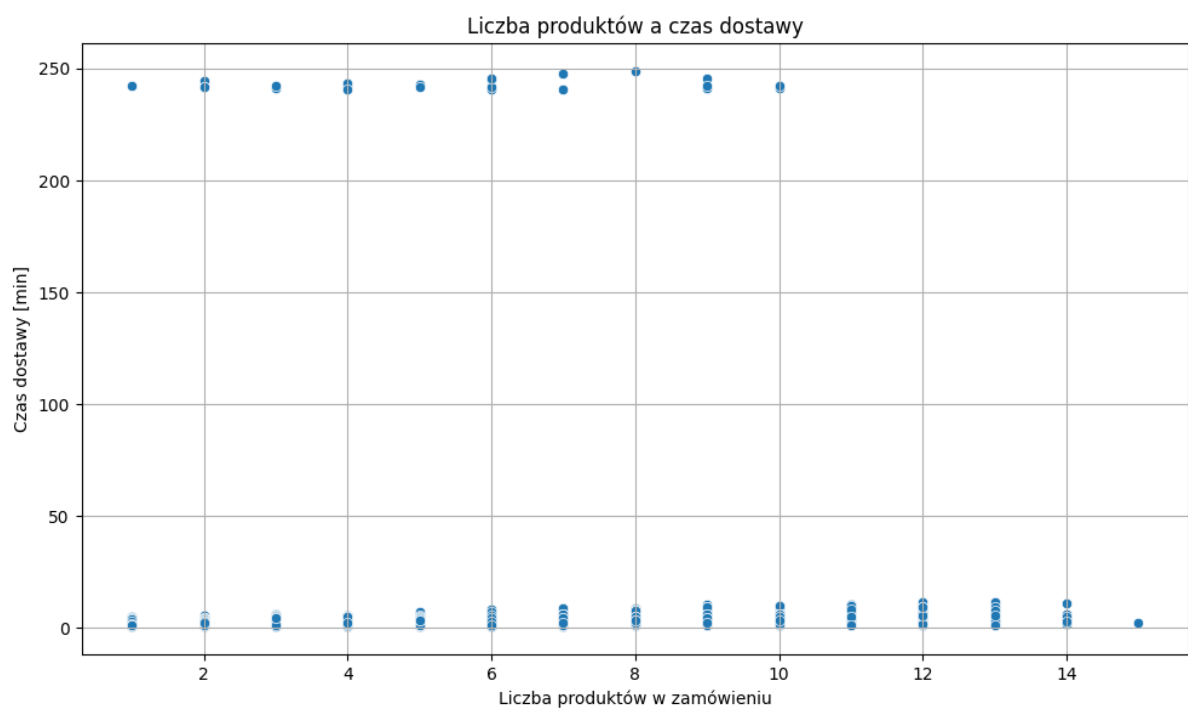
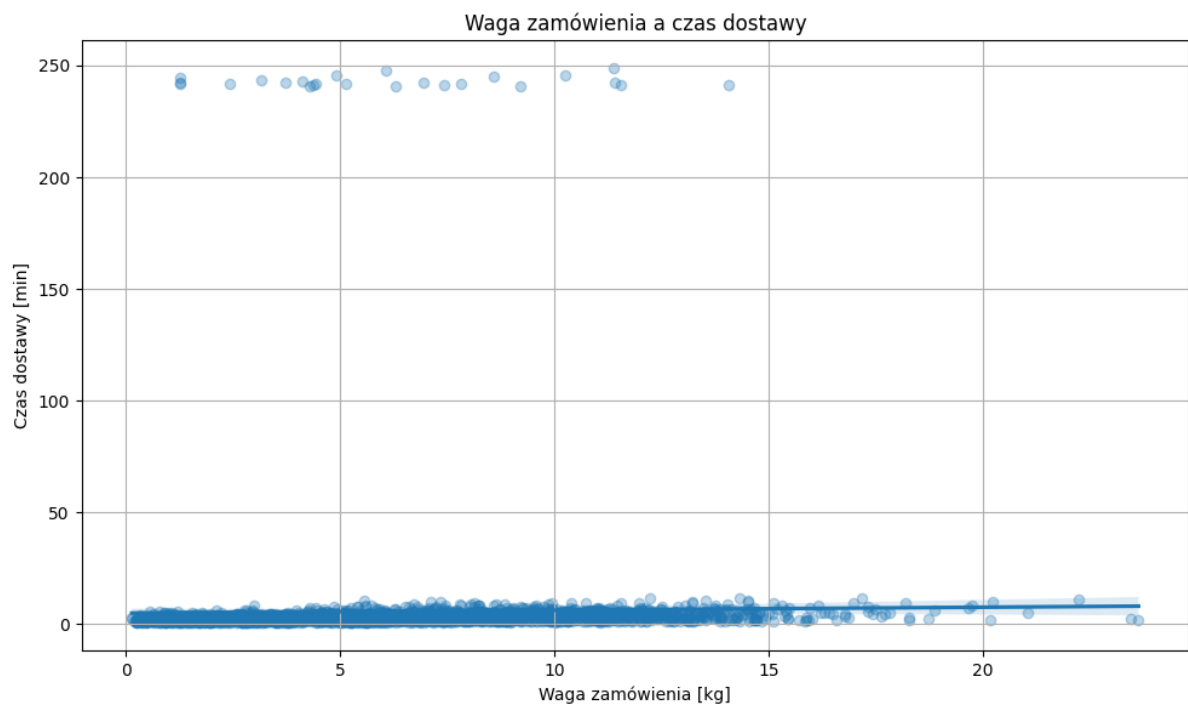
1. **Wczytanie i Czyszczenie Danych:** Importowaliśmy dane dotyczące zamówień, segmentów trasy, produktów oraz relacji między nimi. W procesie czyszczenia usunęliśmy rekordy zawierające braki danych lub błędne wartości.
2. **Filtracja Segmentów „STOP”:** Skupiliśmy się na segmentach, w których nastąpiło zatrzymanie pojazdu. Na ich podstawie obliczono rzeczywisty czas dostawy (na podstawie różnicy między początkiem a końcem segmentu), wyrażony w sekundach, a następnie przekonwertowano i zaokrąglono do minut.
3. **Obliczenie Błędu Predykcji:** Błąd predykcji określaliśmy jako różnicę między działaniem naszego obecnego algorytmu (planowany czas) a rzeczywistym czasem dostawy. Dzięki temu mogliśmy zidentyfikować, że planowany czas (około 3 minut) jest drastycznie niższy od faktycznego (około 240–250 minut).
4. **Agregacja i Wizualizacja:** Przygotowaliśmy szereg wykresów, w tym:
 - Histogram rzeczywistego czasu dostawy (1-minutowa granularity),
 - Histogram błędu predykcji,

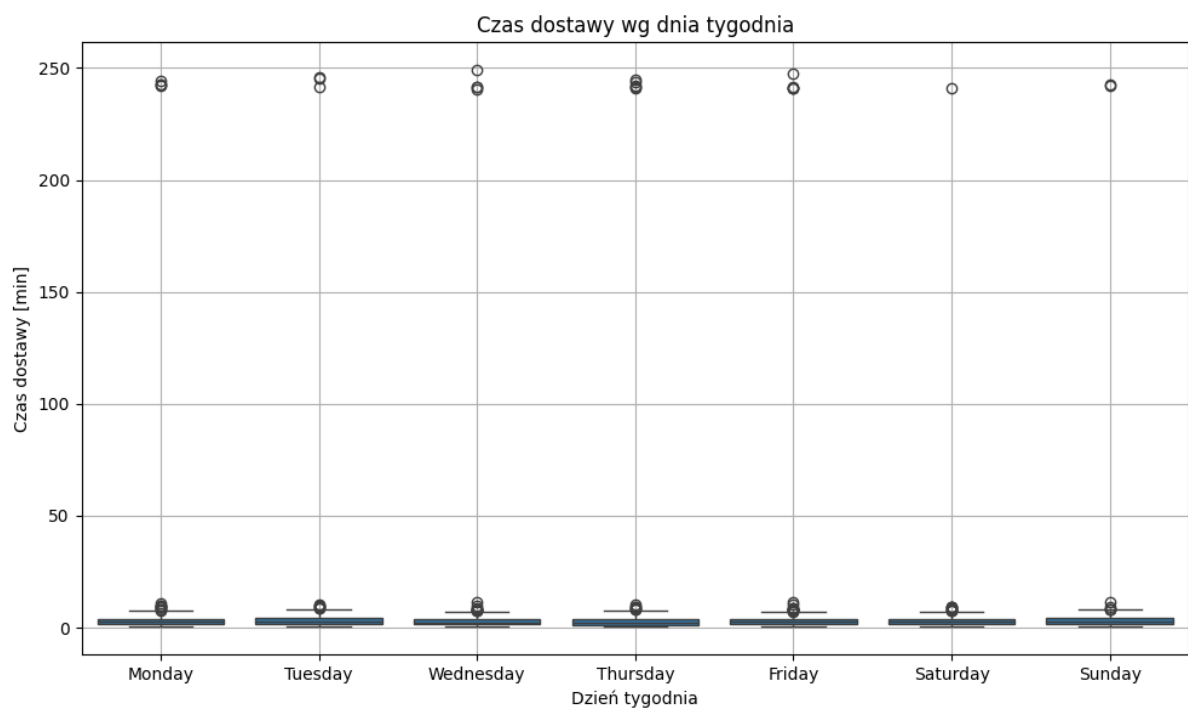
- Wykresy przedstawiające średnie czasy dostawy według sektora, dnia tygodnia i godziny,
- Analizy grupowe zamówień, dla których błąd przewyższa 100 minut.

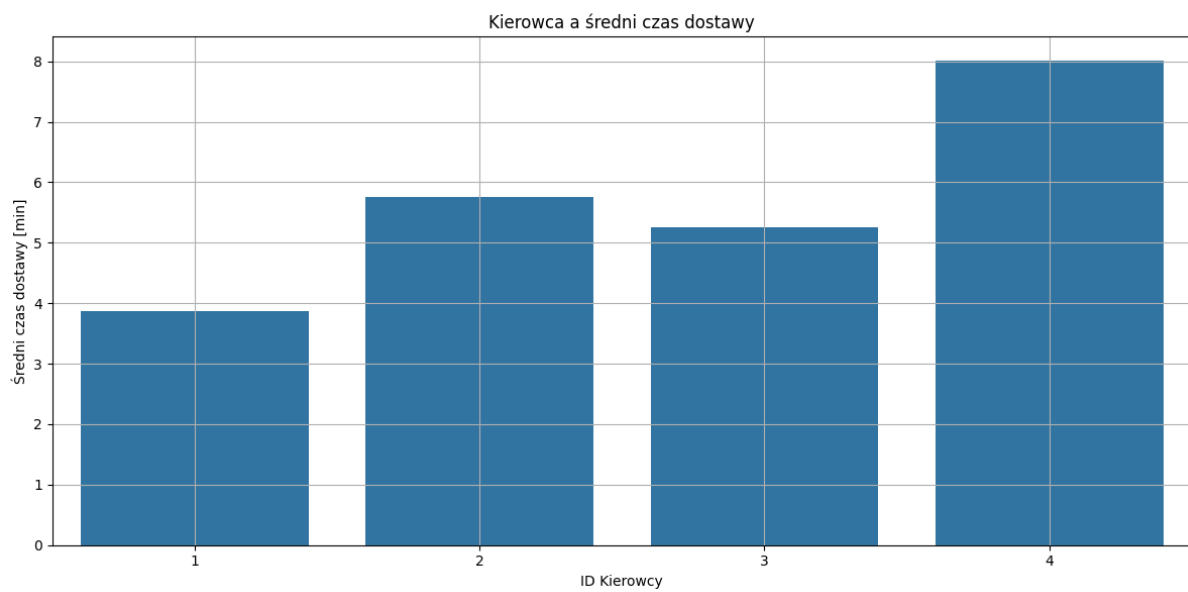
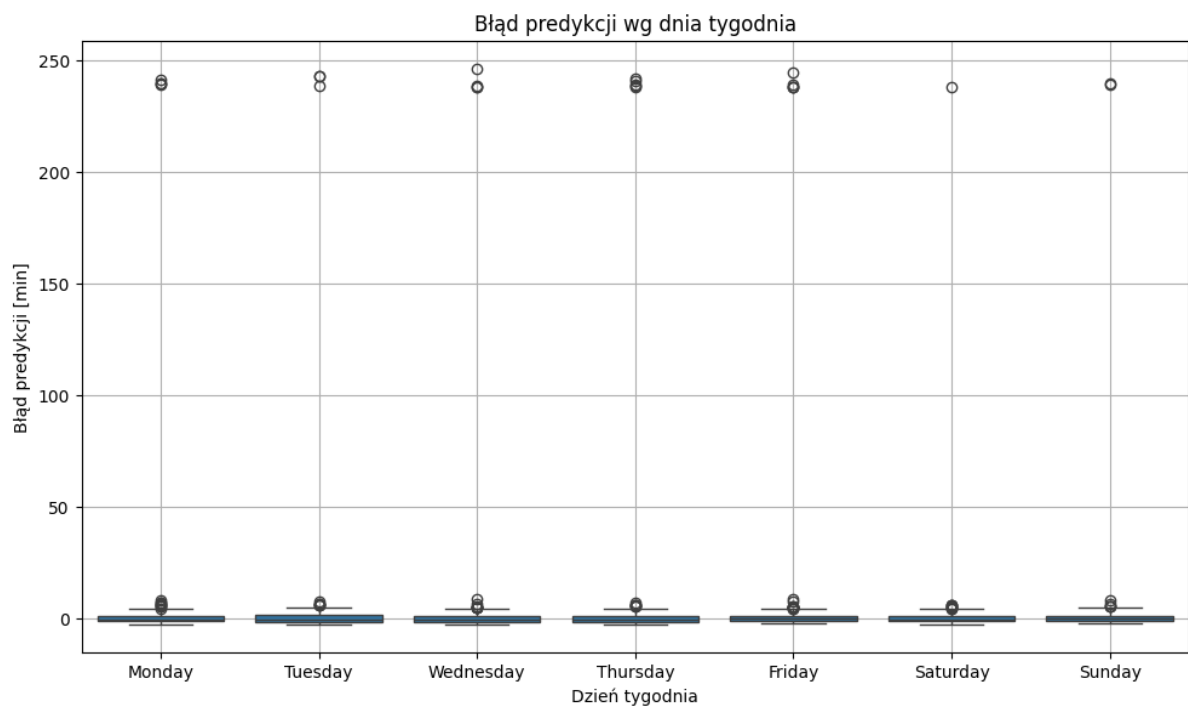
Dzięki tym wizualizacjom mogliśmy wyłapać potencjalne trendy i zależności, które sugerują złożoność procesu dostawy.

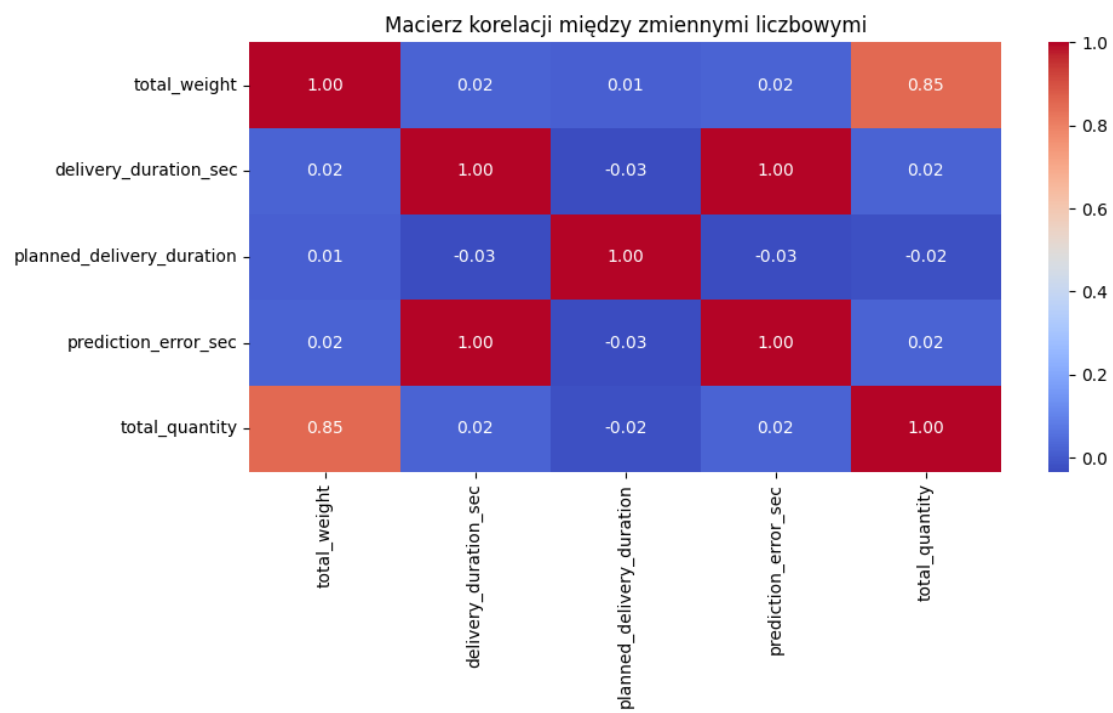
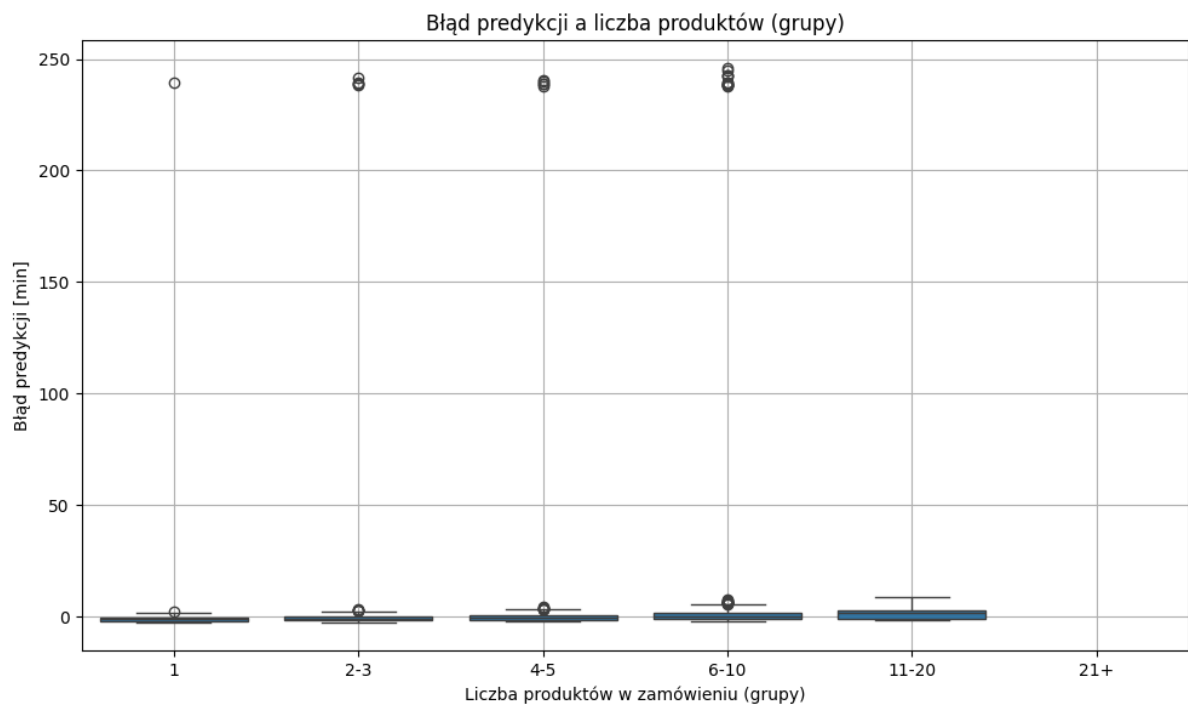


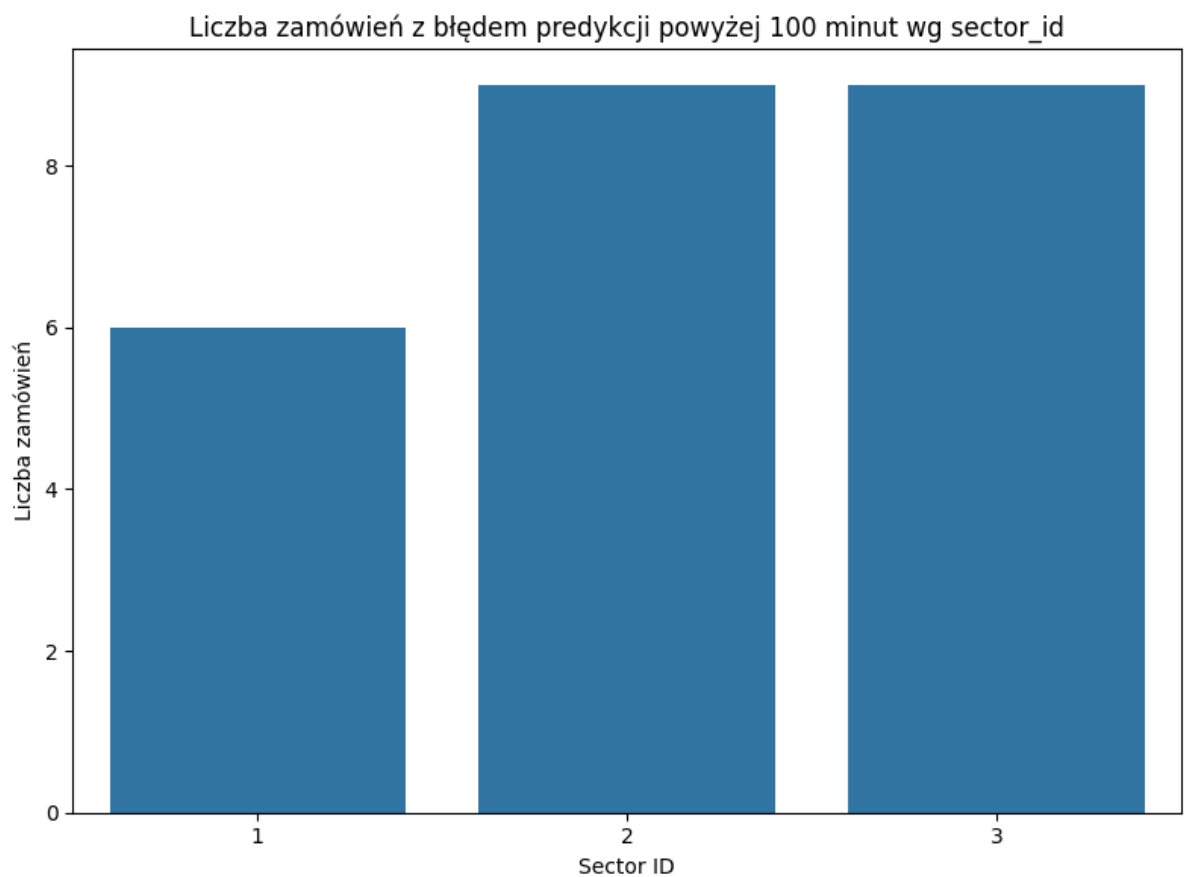
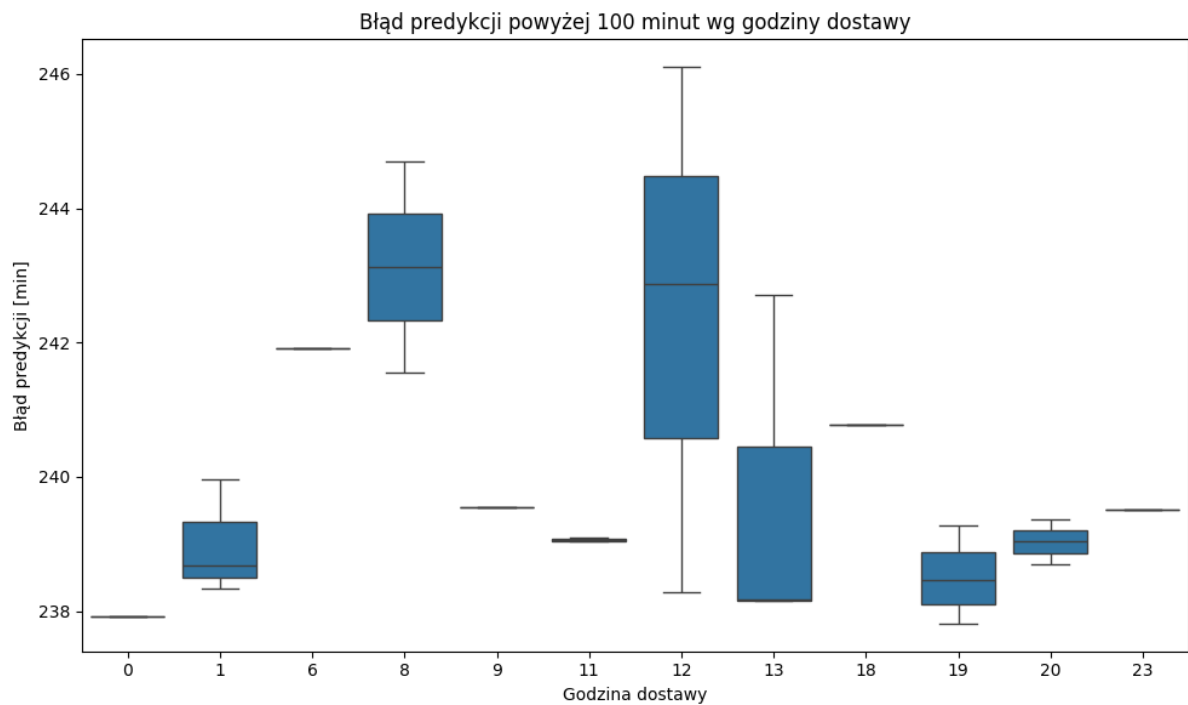












3. Budowanie i Weryfikacja Hipotez

3.1. Predykcja Czasu Dostawy w Oparciu o Sektor

Hipoteza: Prognozowanie czasu dostawy oddzielnie dla każdego sektora może poprawić dokładność, ponieważ różne sektory mogą charakteryzować się różnymi warunkami logistycznymi.

Jak walidować hipotezę:

- **Analiza Statystyczna:** sprawdzamy, czy średnie czasy dostaw znacząco różnią się między sektorami. Na załączonym wykresie „Średni czas dostawy w zależności od sector_id” widzimy, że poszczególne sektory mają wyraźnie różniące się czasy.
- **Testy Statystyczne:** Można zastosować testy istotności (np. test ANOVA), aby potwierdzić statystyczną różnicę między sektorami.
- **Wizualizacja:** Wykresy, na których prezentujemy dane dla poszczególnych sektorów, pozwalają łatwo zidentyfikować obszary o systematycznie dłuższych lub krótszych czasach dostawy.

3.2. Alternatywna Metoda Predykcji

Propozycja: Zamiast stosować prostą średnią, można wykorzystać model regresyjny lub algorytmy uczenia maszynowego, które uwzględniają dodatkowe zmienne, np.:

- Sektor (region dostawy),
- Liczba produktów w zamówieniu,
- Waga zamówienia,
- Godzina realizacji dostawy,
- Dzień tygodnia.

Metodologia walidacji nowego podejścia:

- **Podział danych:** Podziel dane na zbiór treningowy oraz testowy.
- **Trening modelu:** Wytrenuj model regresyjny dokładając więcej zmiennych, które mogą wpływać na czas dostawy.
- **Ocena wydajności:** Porównaj prognozy nowego modelu ze zbiorami testowymi, stosując miary błędu, takie jak średni błąd absolutny (MAE) lub średni błąd kwadratowy (MSE).
- **Porównanie:** Porównaj wyniki nowego modelu z obecnym podejściem opartym na średniej. Jeśli nowy model wykazuje mniejszy błąd, sugeruje to, że dodatkowe zmienne pomagają poprawić dokładność prognoz.

3.3. Czynniki Wydłużające Czas Dostawy

Przykłady:

- **Brak windy:** W budynkach wielorodzinnych, szczególnie tych starszych, brak windy lub obciążone windy mogą znacząco wydłużać czas wejścia do budynku.
- **Rodzaj budynku:** Budynki mieszkalne lub biurowe, gdzie dostawa wiąże się z przechodzeniem przez korytarze, windy lub koniecznością oczekiwania na dostęp, będą miały dłuższy czas realizacji dostawy.

- **Inne czynniki:** Ruch uliczny, zatory drogowe, warunki pogodowe czy skomplikowane trasy mogą również wpływać na czas dostawy.

3.4. Dodatkowe Dane do Zbierania

Propozycje przyszłych zbiorów danych:

- **Typ budynku:** Informacja czy odbiorca mieszka w domu jednorodzinnym, czy w budynku wielorodzinnym.
- **Liczba pięter oraz dostępność windy:** Dane te pozwoliłyby lepiej ocenić, jak logistyka budynku wpływa na czas dostawy.
- **Warunki ruchu drogowego:** Informacje o natężeniu ruchu, zatory, szczególnie w godzinach szczytu.
- **Warunki pogodowe:** Dane o pogodzie, które mogą wpływać na warunki drogowe oraz ogólną dostępność.
- **Dane o trasie:** Informacje dotyczące dokładnej trasy (drogi, skrzyżowania) mogą pomóc w ocenie długości trasy oraz potencjalnych utrudnień.

3.5. Ryzyko Przewidywania (Over- i Under- Estimation)

Ryzyko Niedoszacowania (Under-Estimating):

- Może prowadzić do opóźnień, ponieważ rzeczywisty czas dostawy będzie znacznie przekraczał prognozowany, przez co klienci mogą być zaskoczeni opóźnieniami.
- Negatywnie wpływa to na efektywność pracy kierowców i planowanie zasobów.

Ryzyko Przeszacowania (Over-Estimating):

- Z kolei, przewidywanie zbyt długich czasów dostawy może powodować marnowanie zasobów, kierowcy mogą mieć czas przestoju, a cały system logistyczny może działać nieoptymalnie.
- Może to również wpłynąć na satysfakcję klientów, którzy spodziewają się krótszych czasów realizacji zamówień.

4. Podsumowanie

Nasze badania wykazały, że:

- Użycie jednej globalnej średniej do prognozowania czasu dostawy jest zdecydowanie niewystarczające, prowadzi do ogromnych błędów w przypadku zamówień, gdzie rzeczywisty czas dostawy przekracza 240 minut.
- Predykcja czasu dostawy na podstawie sektorów (oraz potencjalnie innych cech, takich jak liczba produktów, waga zamówienia, godzina realizacji) może znacząco poprawić dokładność prognozy.
- Istnieją liczne czynniki operacyjne (np. brak windy, skomplikowany układ budynku, warunki ruchu), które wpływają na czas dostawy i powinny zostać uwzględnione w bardziej zaawansowanym modelu.

- Zbieranie dodatkowych informacji (np. typ budynku, liczba pięter, dane o ruchu drogowym) w przyszłości pozwoli lepiej zrozumieć i prognozować czas dostawy.

Wdrożenie bardziej zaawansowanego modelu, który wykorzysta powyższe zmienne oraz podzieli zamówienia na wielokryterialne segmenty, powinno przyczynić się do znacznego zmniejszenia ryzyka zarówno niedoszacowania, jak i przeszacowania czasu dostawy, co ostatecznie przełoży się na lepsze planowanie, efektywność oraz zadowolenie klientów.

5. Wnioski

Aby poprawić obecny model prognozowania:

- **Testujemy hipotezę:** o wyznaczaniu prognoz w oparciu o sektor, analizując znaczące różnice między średnimi czasami dostawy w poszczególnych regionach.
- **Proponujemy alternatywne podejście:** oparte na zaawansowanych algorytmach regresyjnych lub modelach uczenia maszynowego, uwzględniających dodatkowe zmienne.
- **Identyfikujemy czynniki:** takie jak brak windy lub skomplikowana infrastruktura budynku, które mogą wydłużać czas dostawy.
- **Rekomendujemy zbieranie dodatkowych danych:** które mogą znacząco wpłynąć na dokładność prognoz.
- **Analizujemy ryzyka:** jakie niesie zarówno niedoszacowanie, jak i przeszacowanie czasu dostawy, aby zoptymalizować planowanie tras oraz zarządzanie zasobami firmy.

Załączone Wykresy:

- **Histogram rzeczywistego czasu dostawy (w minutach):** Pokazuje, że dla niektórych przypadków większość dostaw trwa od 240 do 250 minut.
- **Histogram błędu predykcji (w minutach):** Ilustruje różnicę między planowanym a rzeczywistym czasem, błąd sięga średnio około dla niektórych przypadków 240 minut.
- **Średni czas dostawy wg sector_id:** Wykres słupkowy, który wskazuje na drobne, ale systematyczne różnice między sektorami.
- **Planowany czas vs. błąd predykcji:** Pokazuje, że niska wartość planowanego czasu skutkuje bardzo wysokim błędem.
- **Dodatkowe wykresy:** Analiza korelacji między wagą, liczbą produktów, godziną, a czasem dostawy, oraz zestawienia zbiorcze dla zamówień z bardzo wysokim błędem.

Raport ten prezentuje nasze podejście do budowania i weryfikacji hipotez w kontekście poprawy predykcji czasu dostawy. Proponowane zmiany i zbieranie dodatkowych danych mają na celu osiągnięcie bardziej precyzyjnych prognoz, co wpłynie na optymalizację pracy kierowców oraz poprawę satysfakcji klientów.