

Trash-ML: KI Analyse von Stadtreinigungs Daten zur Optimierung von Einsatzplänen

To Uyen Nguyen Thi | Viet Duc Hoang | Justin Pallas | Krist Baliev

{x} Gliederung

1. Problemstellung
 - 1.1. Problemfrage
 - 1.2. Ziel und Methode
2. Datenrecherche
 - 2.1. Planung der Recherche
 - 2.2. Datenbeschaffung
3. Korrelation zwischen Daten
4. Aggregation der Spalten
5. Modellentwicklung & Evaluierung
6. Fazit und Verbesserungsmöglichkeiten



Berliner Stadtreinigung (Logo)

{1.1} Problemfrage

- Bedarf einer **Optimierung** der Einsatzpläne des **BSR** [Berliner Stadtreinigung]
- Welche **Daten** / **Features** nehmen Einfluss auf die Höhe der **Tonnage** bei der Müllsammlung?



Bild 1: Mülli nicht abgeholt! :(

{1.2} Ziel und Methode

Primäres Ziel

- Bereitstellung eines Prototyp Modells zur Vorhersage von anfallendem Müll für einen jeweiligen Tag
- Training basierend auf internen und externen Daten
- Modelltraining mit bewährten Machine Learning Algorithmen (z.B. Regression)

Methode

1. Planung der Datenrecherche
2. Finden & Aufbereiten von externen Daten
[Recherche, Lag-Einbau, Korrelation]
3. Aggregation von Spalten
4. Modellentwicklung
[Normalisierung, Algorithmuswahl]
5. Evaluierung
6. Eventuelle Modifizierungen

{2.1} Planung der Recherche

Planung der Recherche durch Vermutungen / Hypothesen

- Wetterdaten [Temperatur, Niederschlag, Schneefall]
- Einwohnerzahlen [Wachstum, Wanderung, Geburten/Tode]
- Umsatz im Einzelhandel/Gastgewerbe
- Ferien und Feiertage
- Verbraucherpreisindex
- Tourismus [Übernachtungen, Ankünfte]

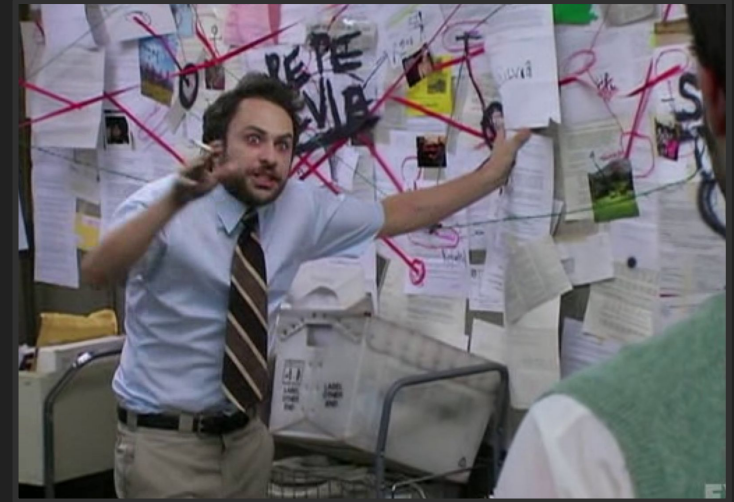


Bild 2: Datenbeschaffung 2025, coloriert

{2.2} Datenbeschaffung

- Statistik Berlin Brandenburg
[Einzelhandel]
- Meteostat Python Bibliothek (Wetterdaten API)
[avg_temperature, precipitation, snow(mm)]
- Genesis / Destatis
[Verbraucherpreisindex, Tourismus]
- Schulferien-online.de
[Ferien, Feiertage]
- Daten.berlin.de
[Retail_Offers]
- Statistische Bibliothek
[Einwohnerzahlen]

```
from datetime import datetime
from meteostat import Daily
import ssl

# Set SSL context.
ssl._create_default_https_context = ssl._create_unverified_context

# Set time period.
start = datetime(2018, 12, 29)
end = datetime(2023, 12, 30)

# Get daily data.
weather_data = Daily('10381', start, end)
weather_data = weather_data.fetch()
```

Bild 3: Wetterdaten durch Meteostat API

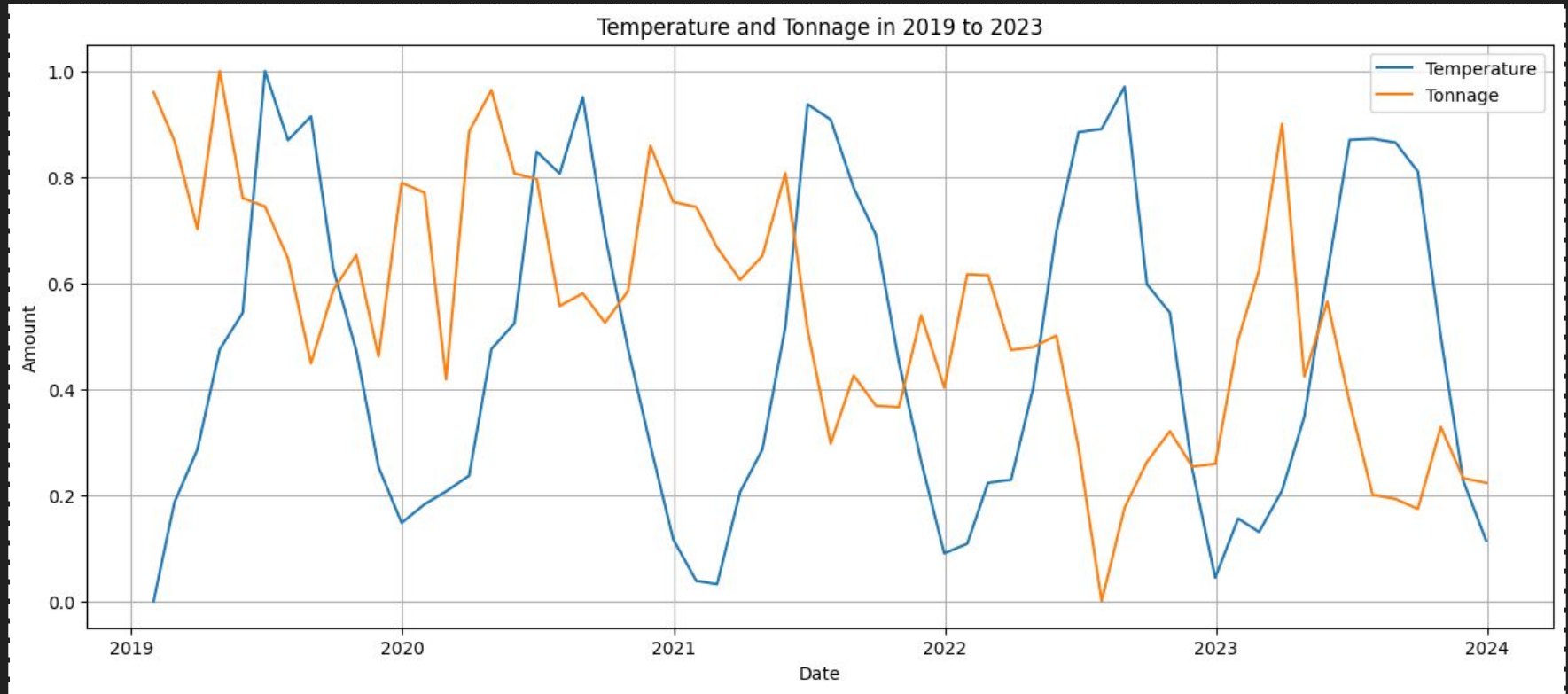
{2.2} Datenaufbereitung

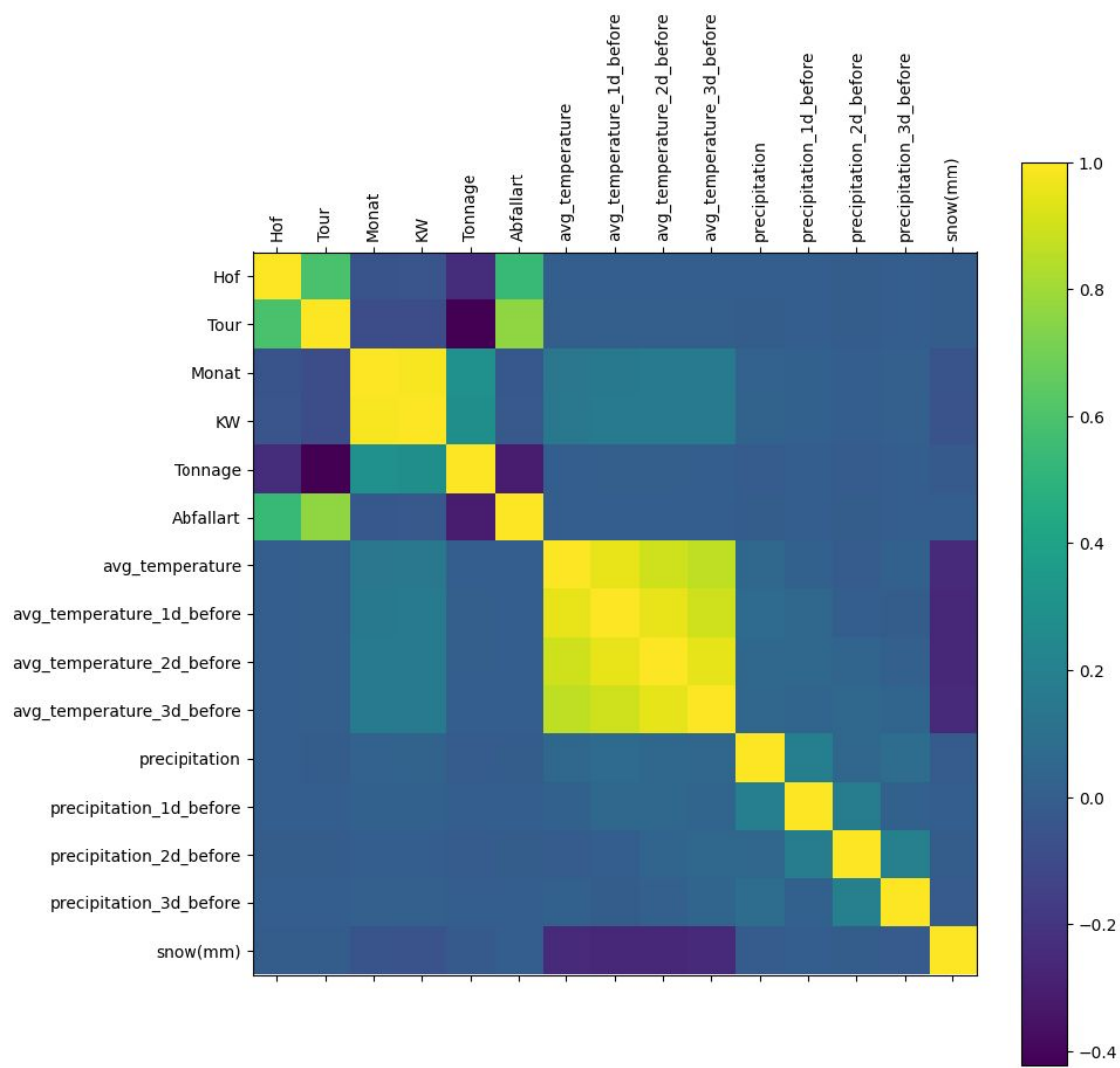
-
- Nan-Values entfernen
- Outlier entfernen
- Datentypen vereinheitlichen
- Mapping auf Spalte Datum

	Month	CW	Year	Yard	Shift	Tour	Tonnage
Date							
13.02.23	602	2107	608923	513	308	44472	5642.73
13.03.23	948	3476	639268	530	323	46916	6805.29
14.02.23	618	2163	625107	529	316	44543	5553.43
14.03.23	960	3520	647360	547	327	44691	7013.56
27.03.23	948	4108	639268	530	323	42526	6091.06
28.03.23	969	4199	653429	550	330	44689	6010.49

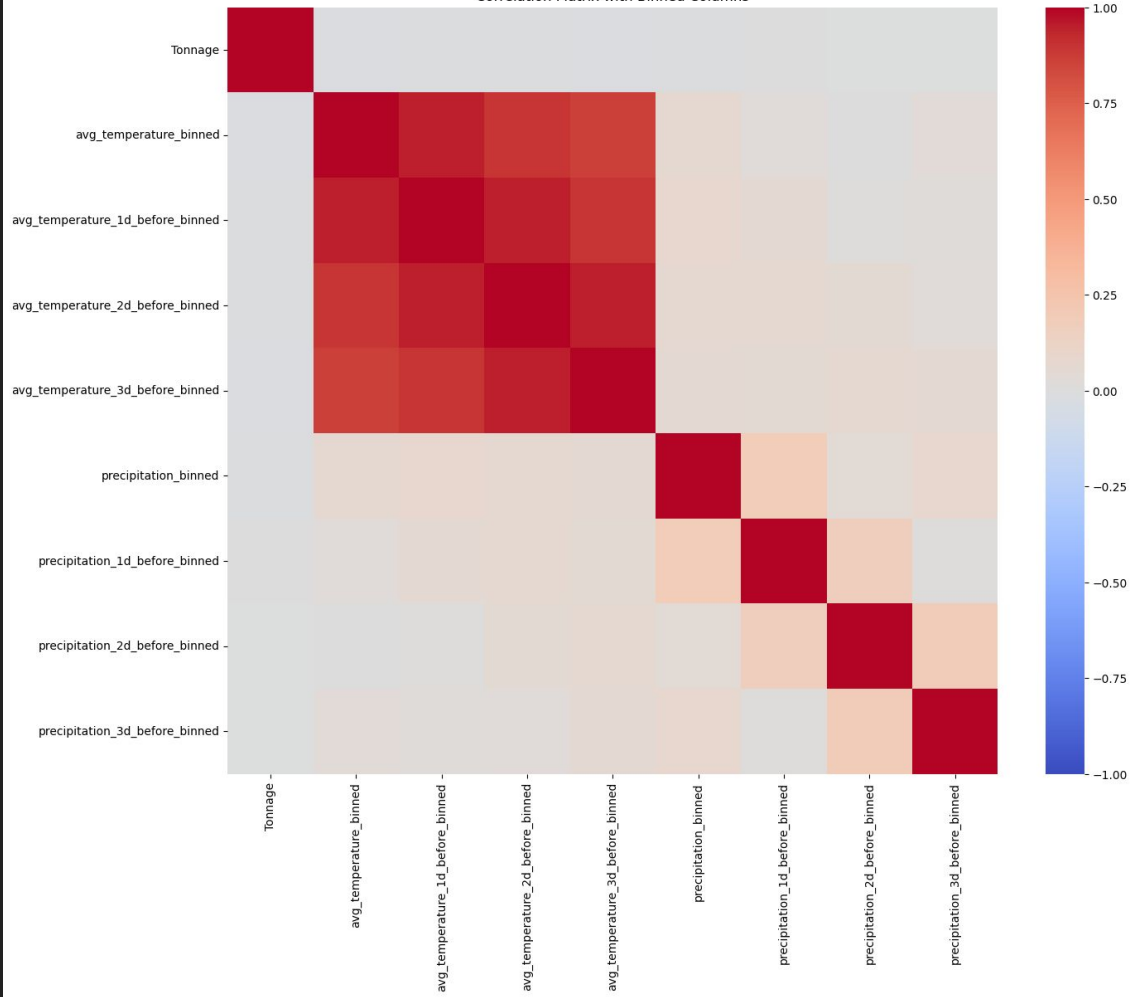
Tage mit einer Tonnage von größer 5500

{3} Korrelation zwischen Daten

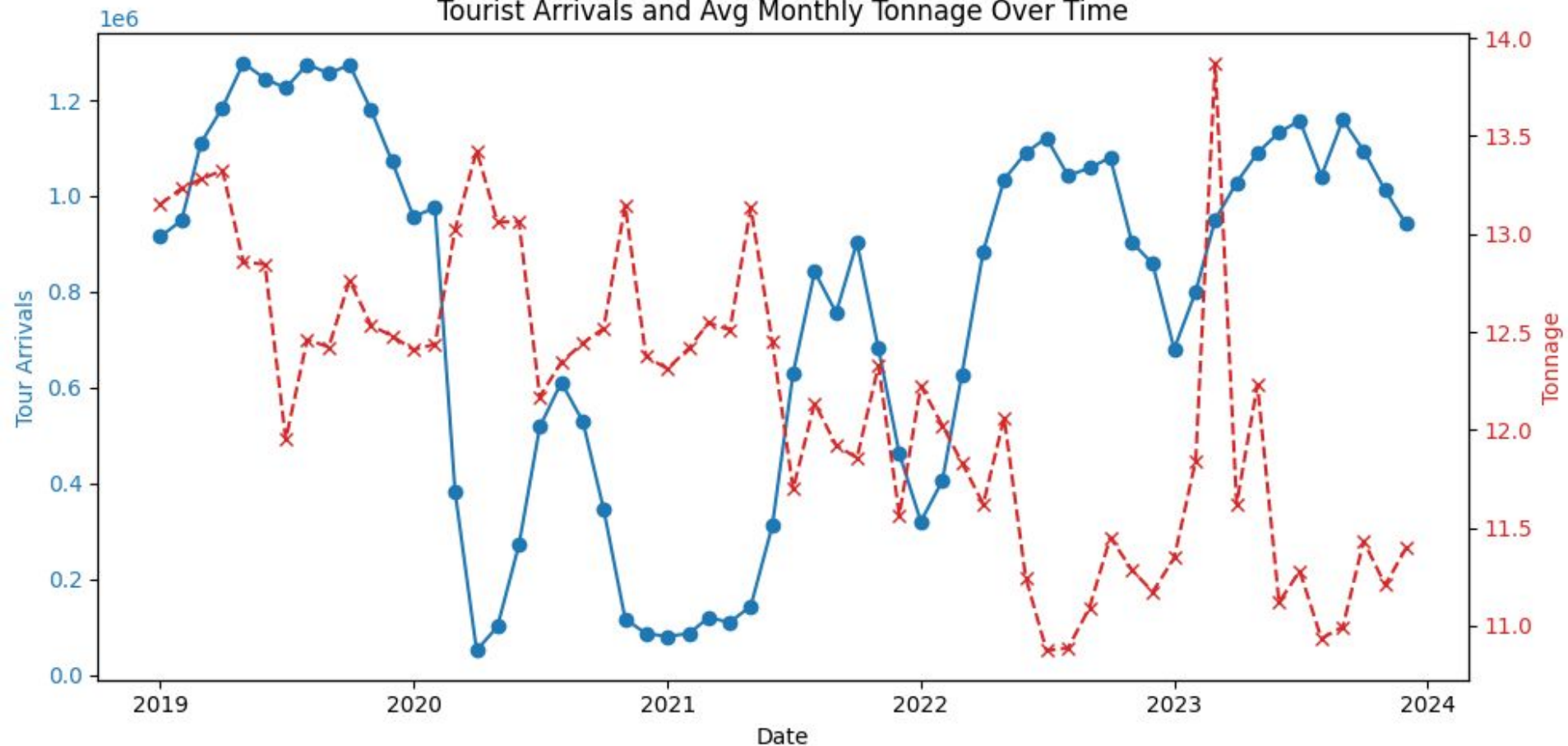




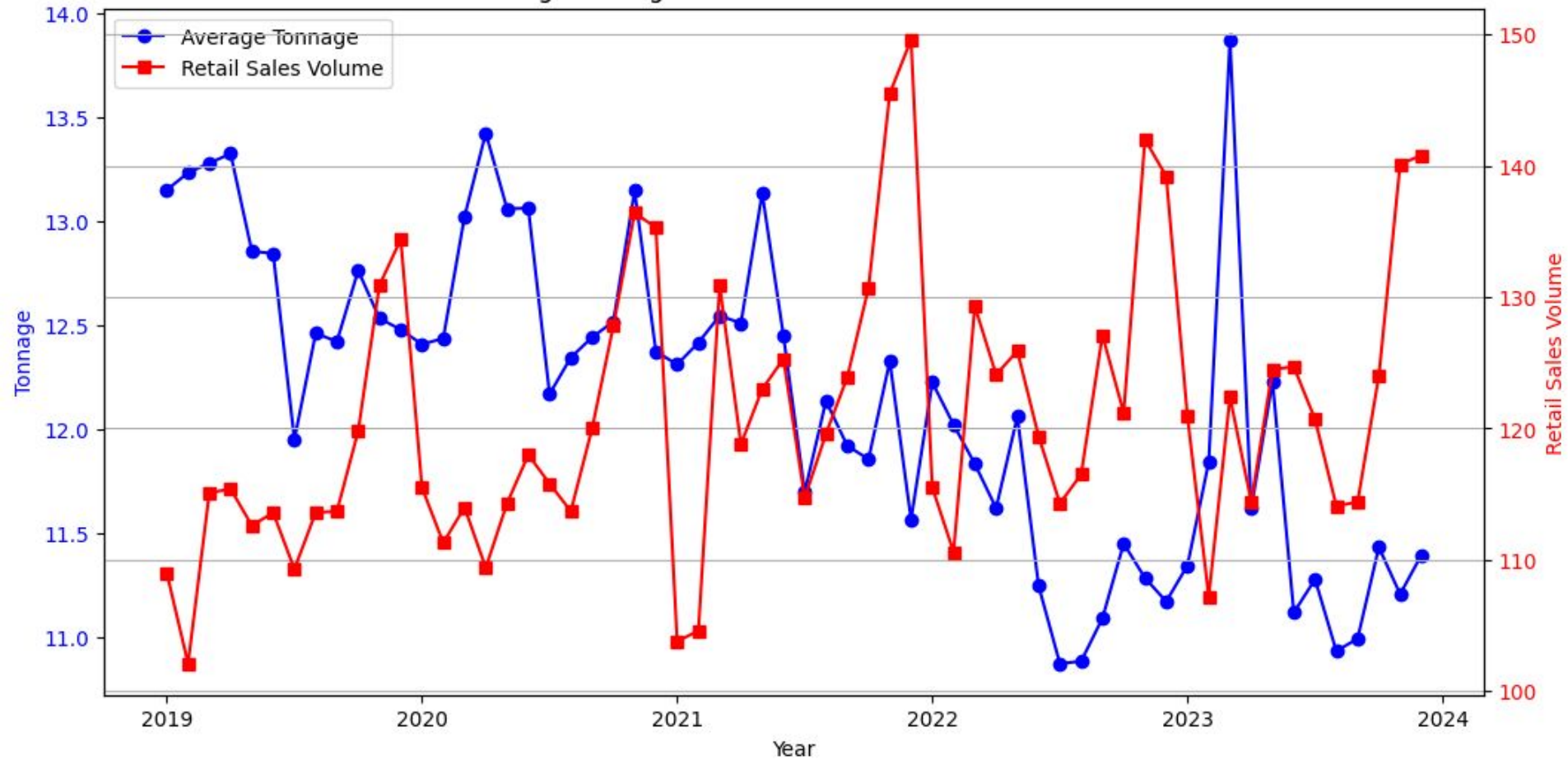
Correlation Matrix with Binned Columns



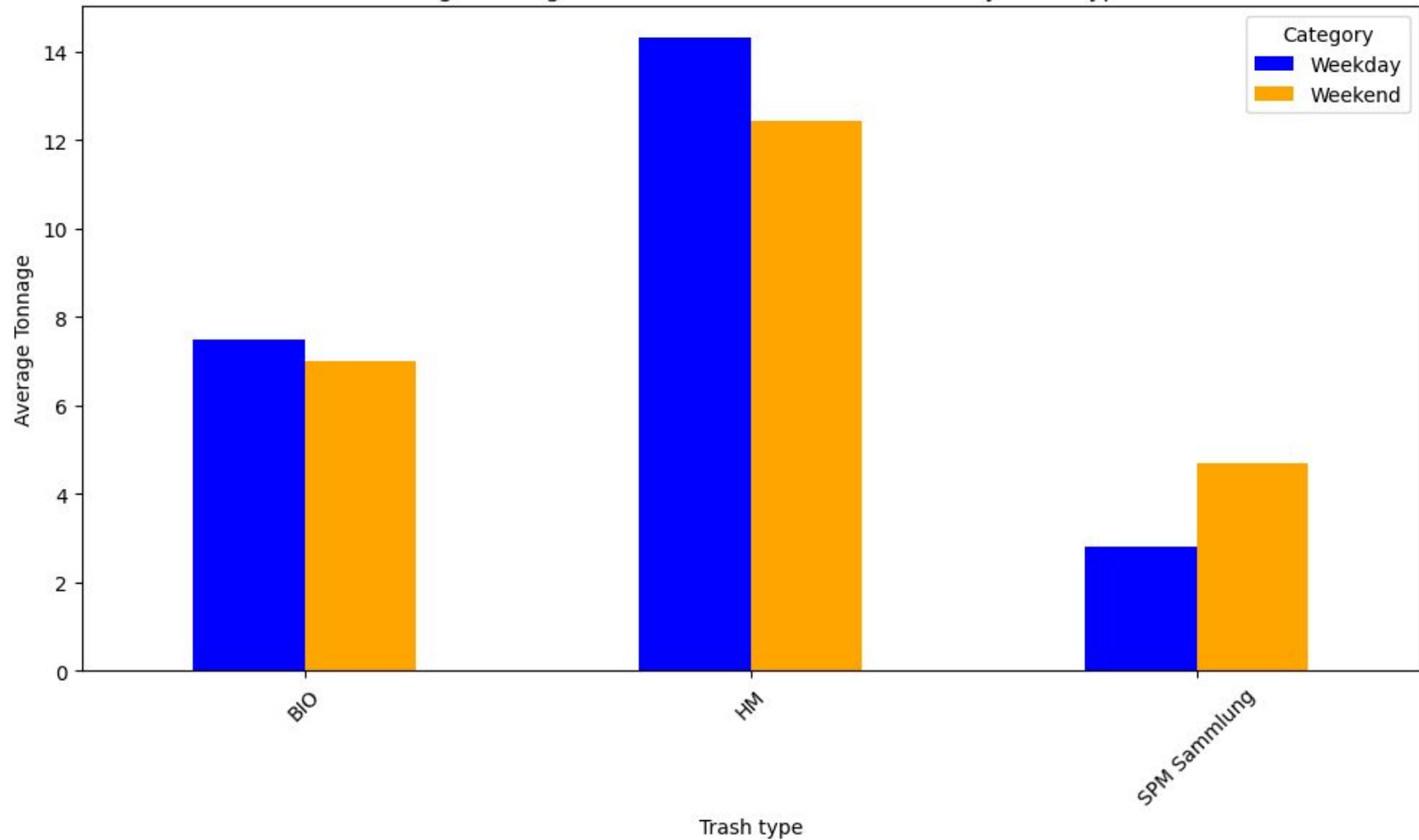
Tourist Arrivals and Avg Monthly Tonnage Over Time



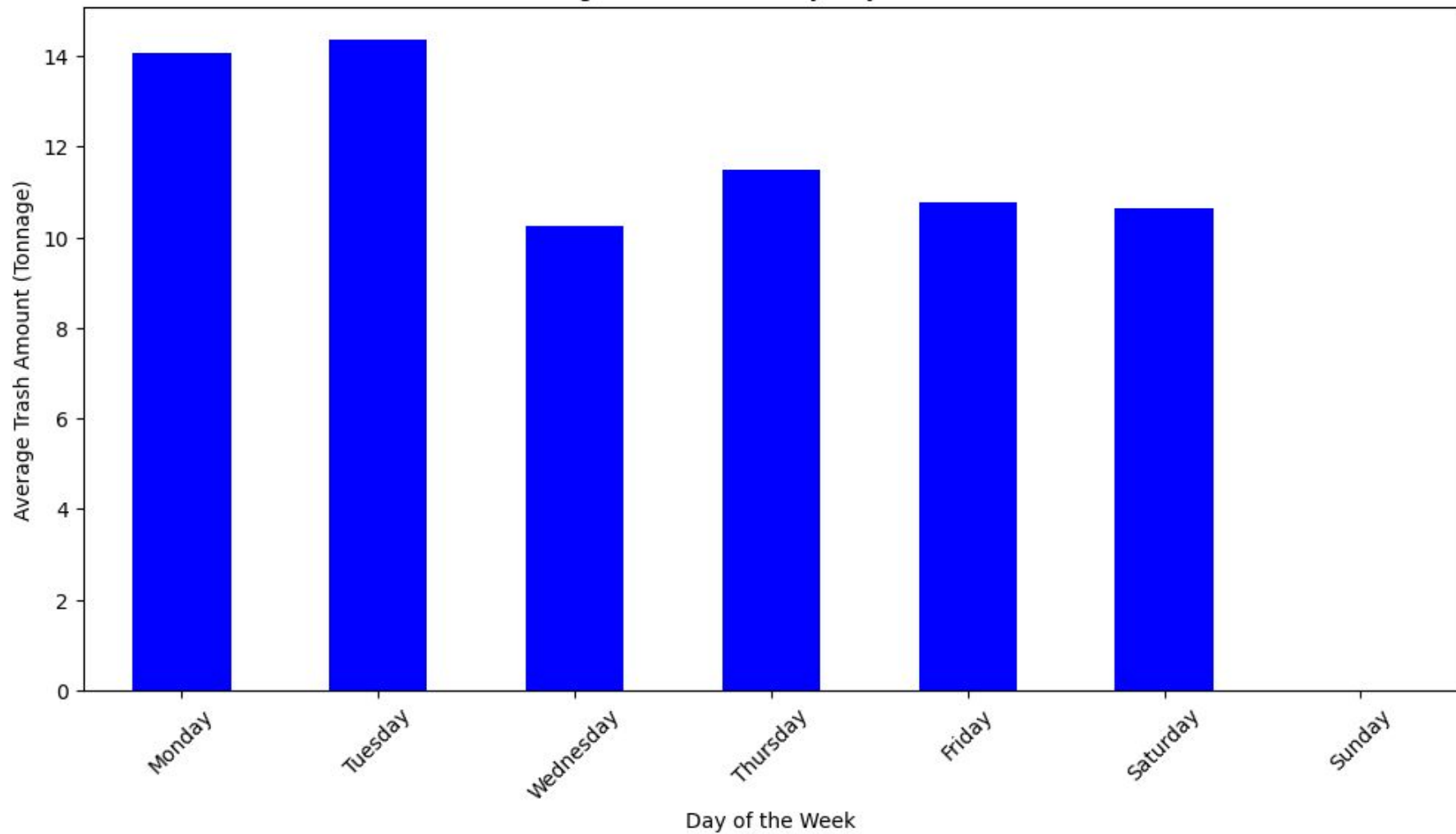
Average Tonnage and Retail Sales Volume Over Time



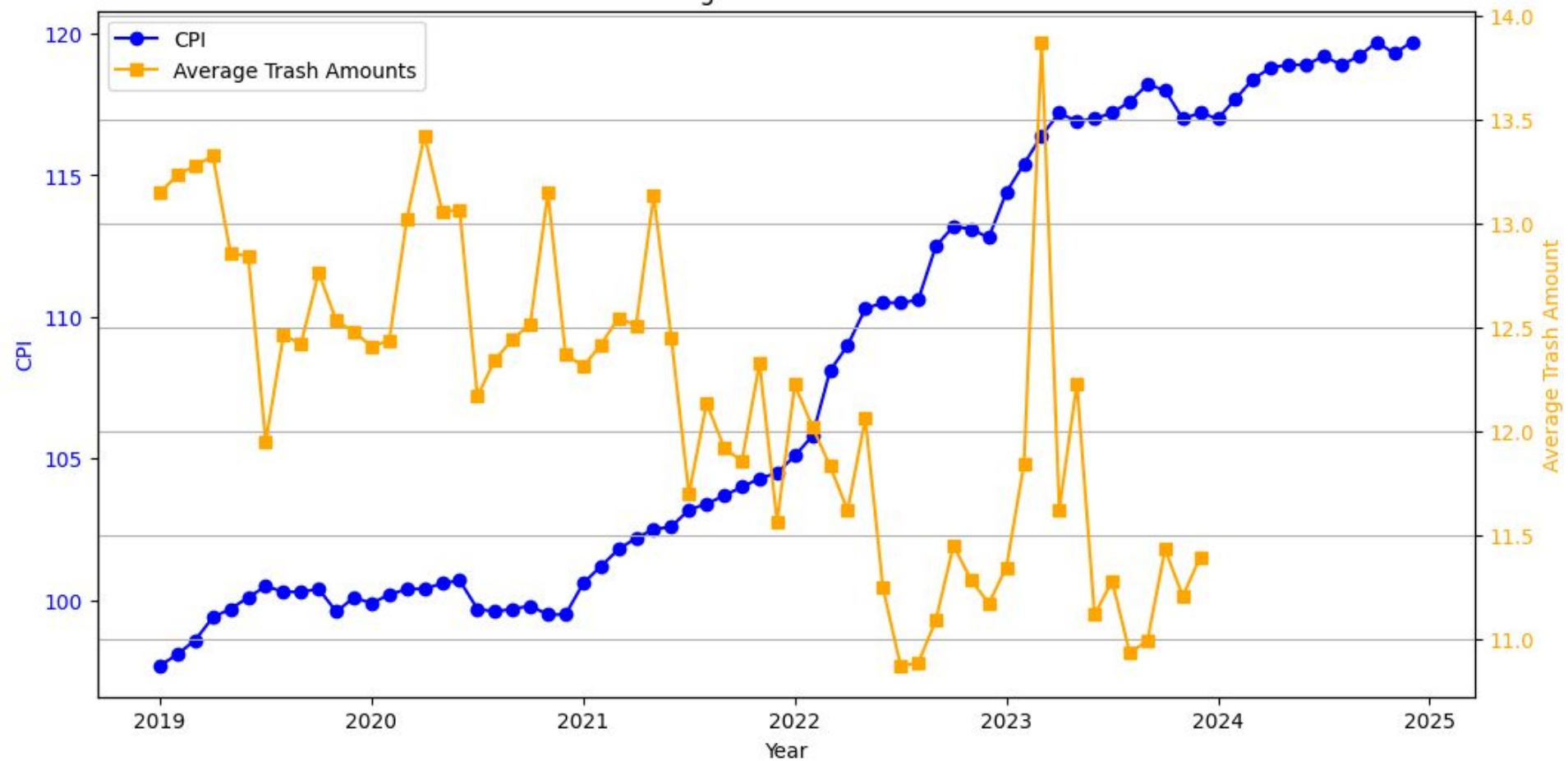
Average Tonnage on Weekends vs Non-Weekends by Trash Type

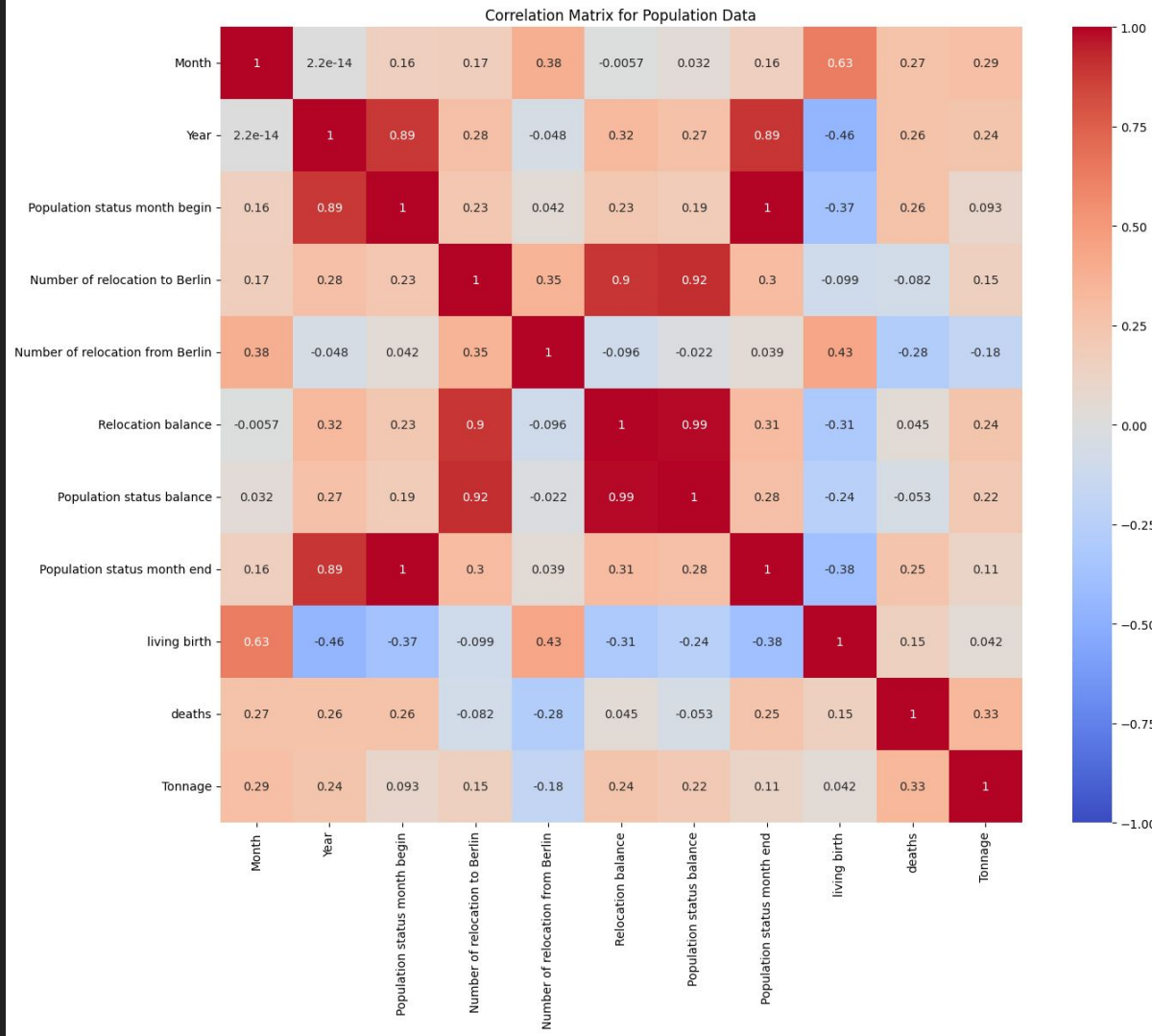


Average Trash Amount by Day of the Week

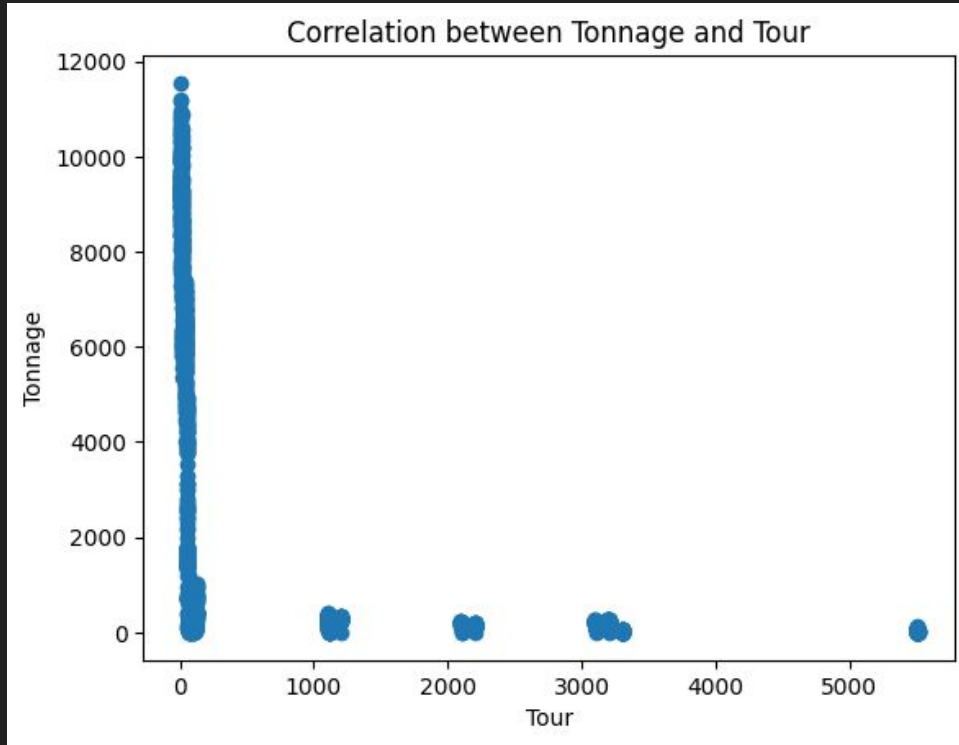


CPI and average trash amounts over time





{3} Auswahl gezielter Features



['SPM Sammlung',
'HM',
'BIO',
'Tonnage',
'Shift',
'Tour',
'Tour_1',
'Tour_2',
'Tour_3',
'Tour_4',
'Tour_5',
'Tour_7',
'Tour_8',
'Tour_9',
'Tour_10',
'Tour_12',
'VMWSN',
'VMWSF',
'deaths',
'weekday_Monday',
'weekday_Tuesday']

{4} Aggregation der Spalten

```
Date,Month,CW,Year,Yard,Shift,Tour,Tonnage,Type,Tour_Stays,avg_temperature,avg_temperature_1d_before,avg_temperature_2d_before,
02.01.19,1,1,2019,VMF,1,1,"5,59",BIO,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,True,Fa
02.01.19,1,1,2019,VMF,1,4,"3,23",BIO,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,True,Fa
02.01.19,1,1,2019,VMF,1,5,"5,68",BIO,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,True,Fa
02.01.19,1,1,2019,VMF,1,6,"5,48",BIO,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,True,Fa
02.01.19,1,1,2019,VMF,1,7,"7,84",BIO,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,True,Fa
02.01.19,1,1,2019,VMF,1,9,"4,63",BIO,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,True,Fa
02.01.19,1,1,2019,VMF,1,10,"4,38",BIO,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,True,Fa
02.01.19,1,1,2019,VMF,1,12,"9,50",BIO,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,True,Fa
02.01.19,1,1,2019,VMF,1,1,"24,00",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,2,"20,70",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,3,"22,30",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,4,"17,16",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,5,"11,60",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,6,"19,10",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,7,"15,94",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,8,"20,94",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,9,"16,94",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,10,"11,28",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,11,"10,54",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
02.01.19,1,1,2019,VMF,1,12,"22,76",HM,2159841.0,1.8,6.4,5.5,6.2,0.1,2.2,0.4,0.1,True,False,False,False,False,False,False,False,Fa
```

Bild 4: Ausschnitt der finalen CSV

{5} Modellentwicklung

1. Aufteilen der Daten in Test-, Trainings- und Validierungsdaten
2. Trainings des Modells mit unterschiedlichen Algorithmen
[Linear, Random_Forest, SGD, DecisionTree]
3. Evaluierung mit geeigneten Metriken
[RMSE, MAE, R^2]
4. Plotten der Ergebnisse
5. Modelloptimierung?

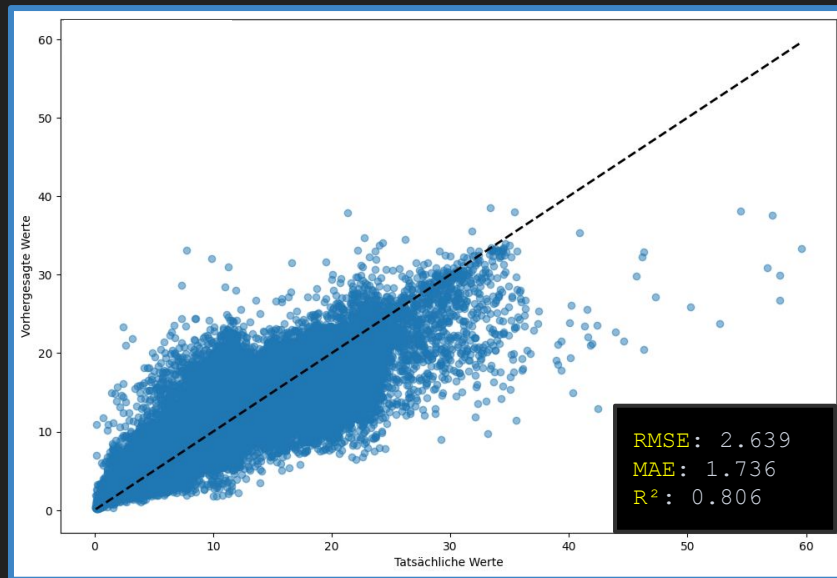
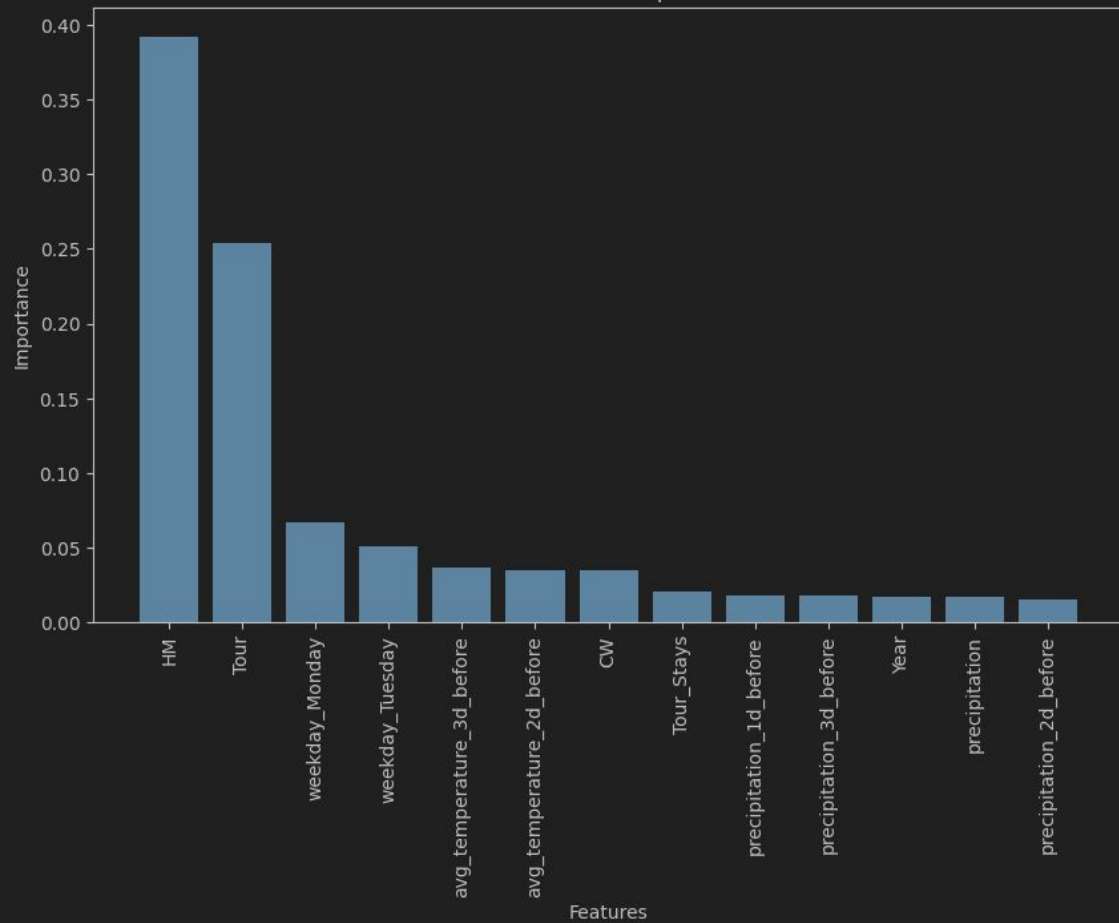


Bild 5: Beste Ergebnisse durch Random Forest Regressor

Filtered Feature Importances



{6} Fazit und Verbesserungsmöglichkeiten

