

Aufgabenblatt III: Datentransformation und Modellierung

1. Vorbereitendes

- a. Erstellen Sie ein neues Verzeichnis mit dem Namen GDSR_3, definieren Sie dieses Verzeichnis als Ihr Arbeitsverzeichnis und speichern dort ein neues R-Skript mit dem Namen `Aufgabenblatt_3.R`.
- b. Ergänzen Sie in Ihrem R-Skript folgende erste Zeile: `rm(list=ls())`
- c. Installieren Sie das Paket `nycflights13` aus dem wir die Daten nutzen werden.
- d. Laden Sie die Pakete `tidyverse` und `nycflights13`.
- e. Betrachten Sie den über `nycflights13` geladenen Datensatz `flights`.
- f. Informieren Sie sich über den Datensatz mit dem Befehl `?flights`. Der Datensatz enthält Informationen zur Pünktlichkeit von allen Flügen, die 2013 von einem der drei New York City-Flughäfen gestartet sind.
- g. Der Datensatz `flights` enthält knapp 10.000 Beobachtungen, in denen eine oder mehrere Informationen fehlen. Erzeugen Sie ein neues R-Objekt mit dem Namen `Fluege_df`, in dem Sie alle vollständigen Beobachtungen des Datensatzes `flights` ablegen. Arbeiten Sie im Folgenden mit dem `Fluege_df` Datensatz weiter.

2. Analyse

Nutzen Sie, wenn möglich, für alle Aufgaben den Pipe-Operator:

- a. Erstellen Sie eine Tabelle mit der Anzahl an Flügen je NYC-Flughafen (`origin`), die eine Verspätung bei Abflug (`dep_delay`) von über 120 Minuten haben.
- b. Erstellen Sie eine Tabelle mit den folgenden Informationen je NYC-Flughafen:
Durchschnittliche Verspätung bei Abflug, minimale Verspätung bei Abflug, maximale Verspätung bei Abflug, Summe der Verspätungsminuten bei Abflug, Anzahl der Flüge.
- c. Erstellen Sie drei Streudiagramme (getrennt nach Flughäfen), die den Zusammenhang zwischen Verspätung bei Abflug und Verspätung am Ziel (`arr_delay`) darstellen.
- d. Erstellen Sie die gleichen Abbildungen noch einmal, aber jetzt nur für Flüge mit mehr als 60 min Verspätung bei Abflug.
- e. Schätzen Sie ein lineares Regressionsmodell mit der abhängigen Variable `arr_delay` und den erklärenden Variablen `dep_delay` und `air_time` und lassen Sie sich den detaillierten Output anzeigen.
- f. Schätzen Sie das Modell aus Teilaufgabe e) noch einmal getrennt für jeden der drei Flughäfen und nur für Flüge mit mehr als 120 min Verspätung bei Abflug.
- g. Speichern Sie das R-Skript ab und prüfen Sie, ob Ihre Ergebnisse replizierbar sind.