

ETL-Prozess anhand von Inlandflügen der USA (2015)

CIP-Projekt Gruppe 16

Hochschule Luzern – Wirtschaft

MSC Applied Information and Data
Science

Data Collection, Integration,
Preprocessing – Projektarbeit (LN2)

Agenda



**Aufgaben-/
Fragestellung**



Datenquellen



Tools



ETL-Prozess

Python
Tableau Prep
SQL Server



**Analyse mit Tableau
Desktop**



Fazit und Reflexion

Aufgaben-/ Fragestellung

Aufgabenstellung:

Generieren eines **ETL**-Prozesses:

- **E**xtract:
 - Datenbeschaffung aus mind. zwei Quellen
- **T**ransform:
 - Aufbereitung und Bereinigung der Daten
 - Verknüpfung der Datensätze
- **L**oad:
 - Daten in Datenbank importieren

Fragestellung:

Welche der Airlines haben die meisten pünktlichsten, verspäteten, annullierten und umgeleiteten Flüge verbucht.

Datenquellen

- Inlandflüge 2015 der USA
 - flights.csv
 - airlines.csv
 - airports.csv
- Rufzeichen der Airlines
 - Gecrawlt mit Python

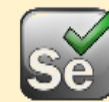
kaggle



Software / Verwendete Tools

Extract

- Python – Selenium – BeautifulSoup



Transform

- Python – Pandas - Numpy
- Tableau Prep



Load

- Microsoft SQL-Server

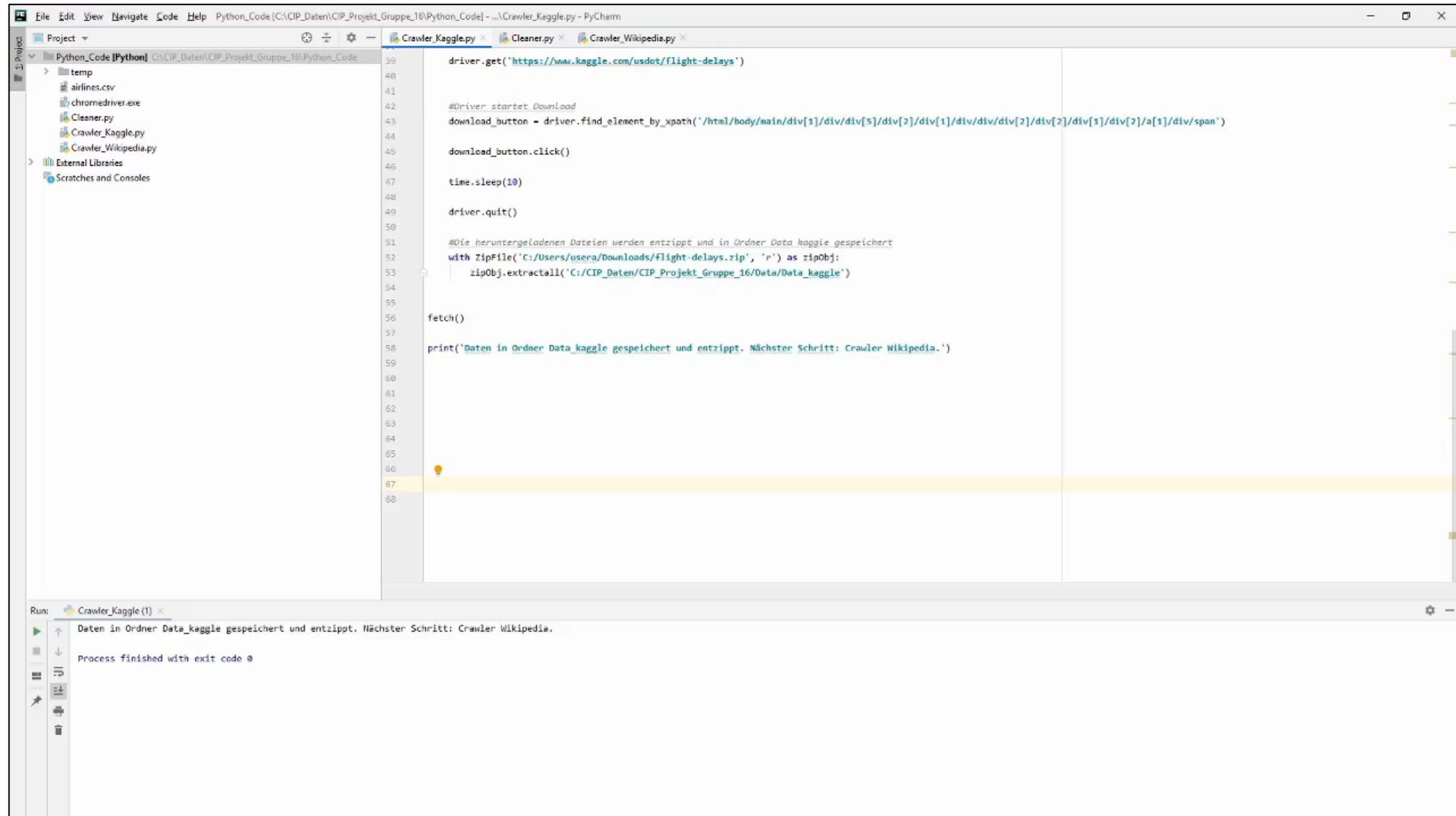


Visualise

- Tableau Desktop



Crawler



Datenbereinigung mit Python

- Analyse der Daten
- Aufbereitung und Bereinigung der Daten

Herausforderungen:

- Datenmenge
- Unvollständige Datensätze
- Nicht eindeutig zuweisbare Einträge

```
In [4]: print(len(df_main))
```

```
5819079
```

```
In [5]: df.head()
```

```
Out[5]:
```

	YEAR	MONTH	DAY	DAY_OF_WEEK	AIRLINE	FLIGHT_NUMBER	TAIL_NUMBER	ORIGIN_AIRPORT	DESTINATION_AIRPORT	SCHEDULED
0	2015	1	1	4	AS	98	N407AS	ANC	SEA	
1	2015	1	1	4	AA	2336	N3KUAA	LAX	PBI	
2	2015	1	1	4	US	840	N171US	SFO	CLT	
3	2015	1	1	4	AA	258	N3HYAA	LAX	MIA	
4	2015	1	1	4	AS	135	N527AS	SEA	ANC	

```
5 rows × 11 columns
```

```
In [7]: # TESTING
is_october = df['MONTH'] == 10

a = df[is_october]
a.head()
```

```
Out[7]:
```

	YEAR	MONTH	DAY	DAY_OF_WEEK	AIRLINE	FLIGHT_NUMBER	TAIL_NUMBER	ORIGIN_AIRPORT	DESTINATION_AIRPORT	SCHEDULED
4385712	2015	10	1	4	AA	1230	N3DBAA	14747	11298	
4385713	2015	10	1	4	DL	1805	N696DL	14771	13487	
4385714	2015	10	1	4	AA	840	N171US	14771	13487	

Transform

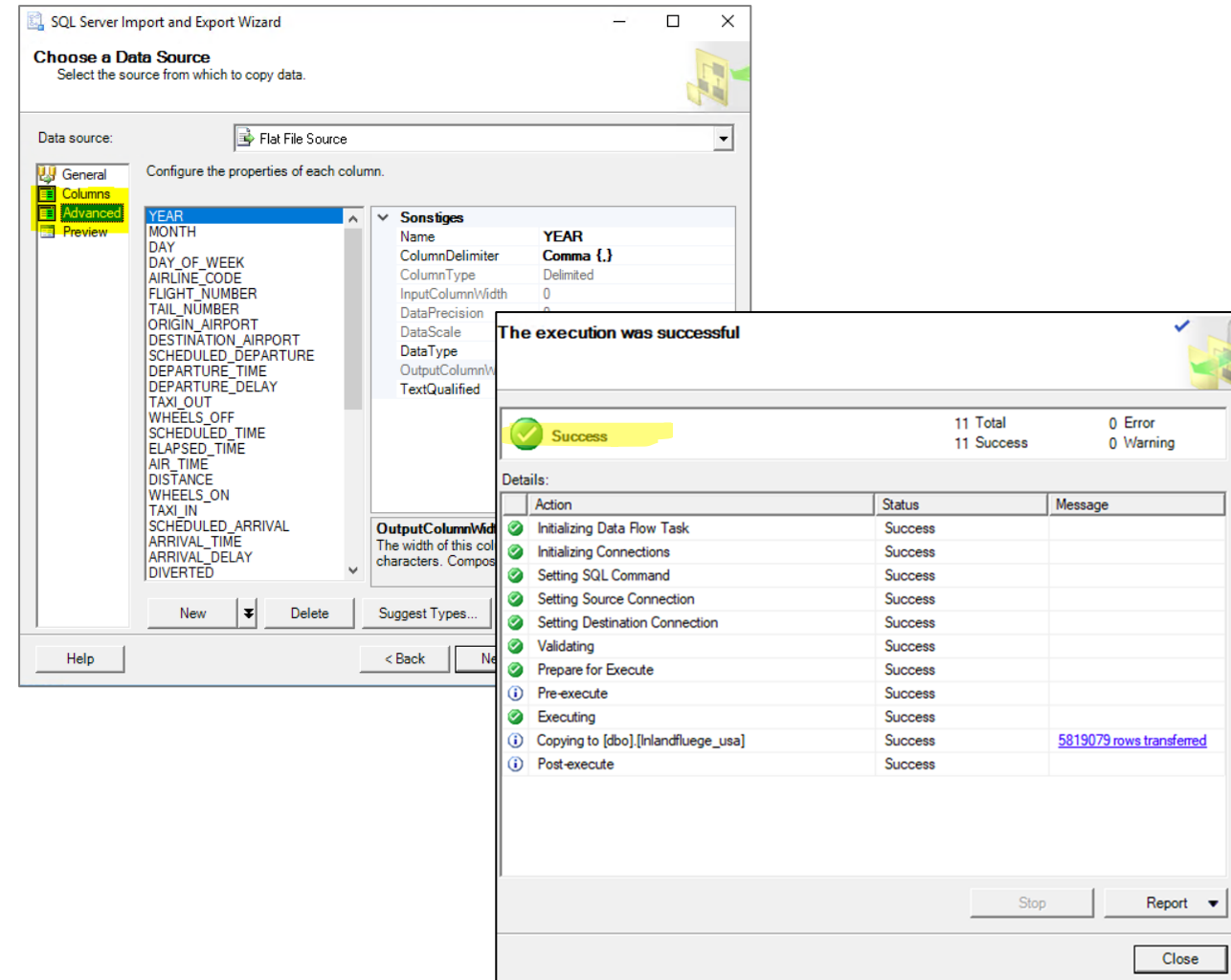
- Zusammenführen der Dateien und Anpassen der Spalten zu einer Ausgabedatei
- Herausforderungen:
 - Die richtige Auswahl der Verknüpfung (left, inner oder right join)
 - Anpassung der Spalten nach jedem join

The screenshot displays the Tableau Prep Builder interface. The top pane shows a workflow: 'flights_new' and 'airports' are prepared (Aufbereiten 1, 2) and then joined (Verknüpfen 1) using an inner join. This result is then joined with 'airlines_merged' (Verknüpfen 2) using a left join. The final result is prepared (Aufbereiten 3, 4, 5) and output (Ausgabe). The bottom pane shows the output configuration, where the data is saved as 'Inlandfluege_usa.csv'. A preview of the CSV data is shown below.

YEAR	MONTH	DAY	DAY_OF_WEEK	AIRLINE_CODE	FLIGHT_NUMBER	TAIL_NUMBER	SCHEDULED_DEPARTURE	DEPARTURE_TIME	DEPARTURE_DELAY
2,015	1	5	1	DL	2,594	N370NB	815	815	0
2,015	1	5	1	DL	2,619	N993DL	815	809	-6
2,015	1	5	1	DL	1,670	N969DL	815	814	-1
2,015	1	5	1	DL	1,692	N802DN	815	816	1
2,015	1	5	1	DL	1,966	N554NW	815	818	3
2,015	1	5	1	DL	2,120	N994AT	815	815	0
2,015	1	5	1	DL	2,275	N3772H	815	816	1
2,015	1	5	1	DL	1,184	N974DL	815	808	-7
2,015	1	5	1	DL	1,497	N6704Z	815	811	-4
2,015	1	5	1	EV	3,262	N14542	815	835	20
2,015	1	5	1	EV	4,141	N17196	815	858	43
2,015	1	5	1	EV	4,737	N21144	815	815	0
2,015	1	5	1	EV	5,077	N851AS	815	809	-6
2,015	1	5	1	EV	5,167	N852AS	815	812	-3

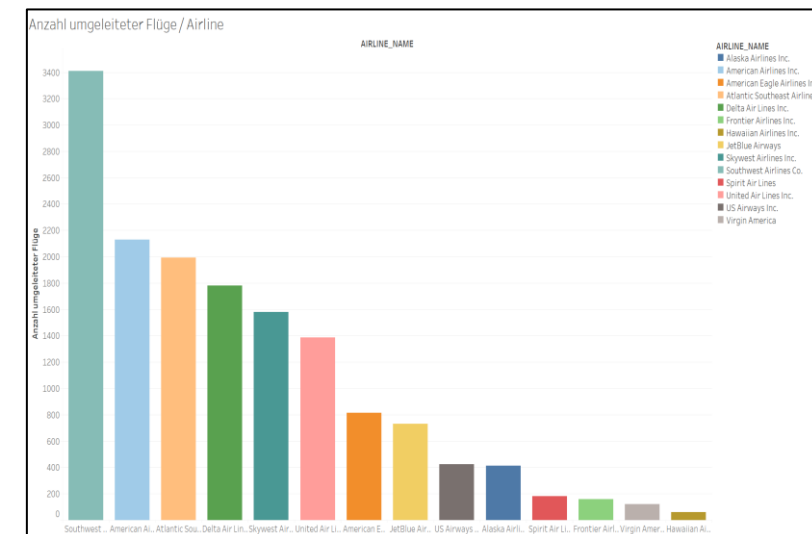
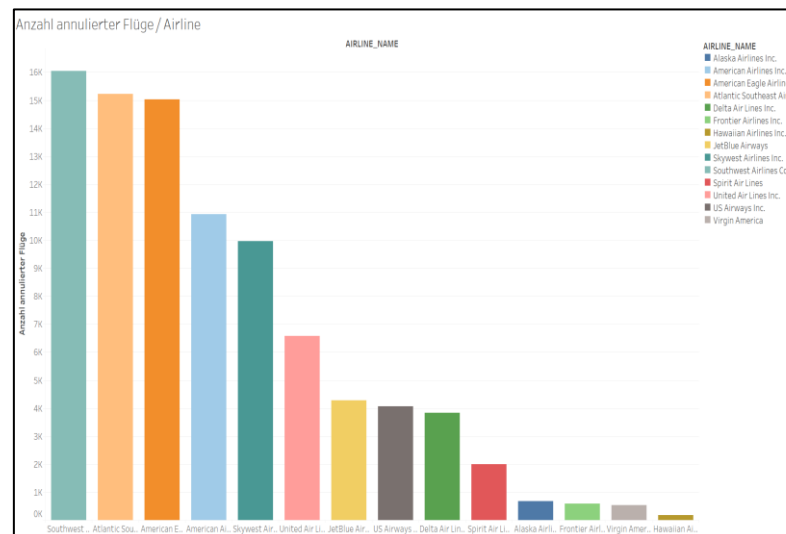
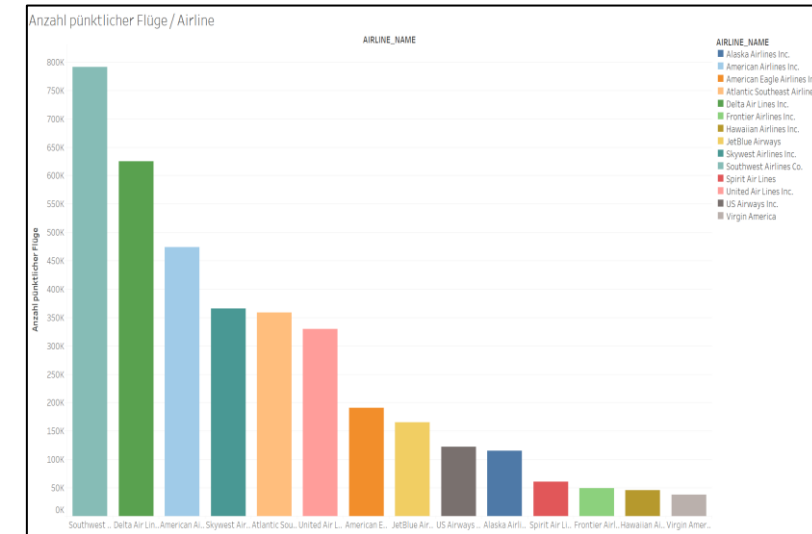
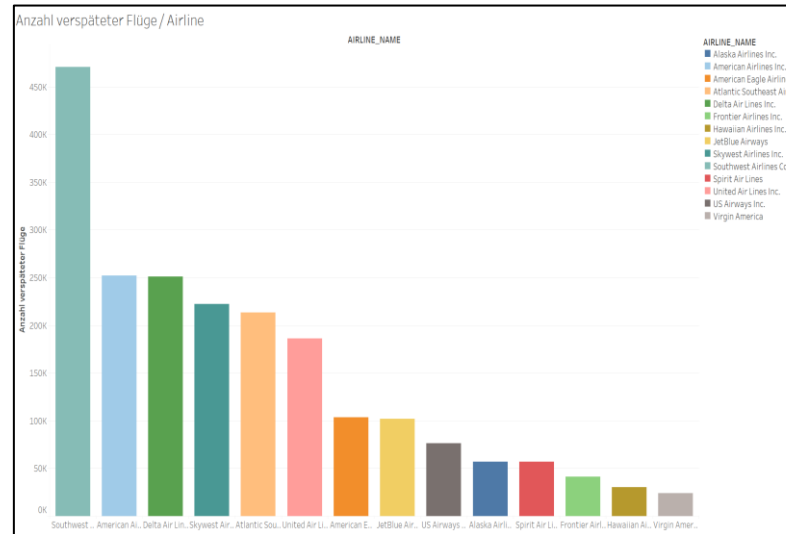
Load | SQL-Server Import

- Erstellung einer Datenbank CIP:
- Import des CSV-Files via Wizard
- Herausforderungen:
 - Wenige Kenntnisse
 - Anpassung der Datentypen und Feldlänge
 - Dauer des Ladeprozesses
- Erfolgreicher Import



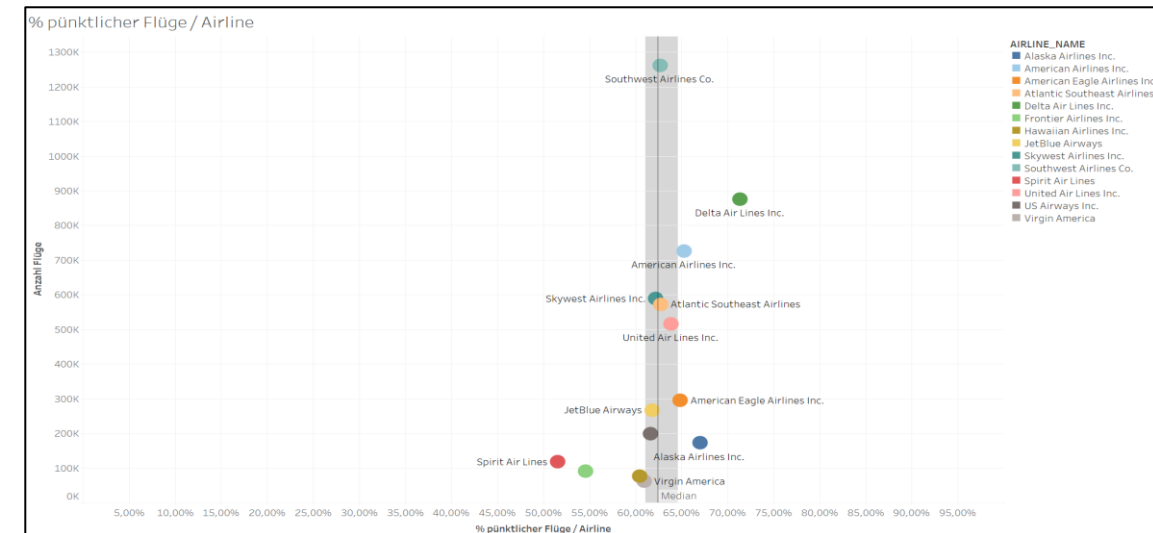
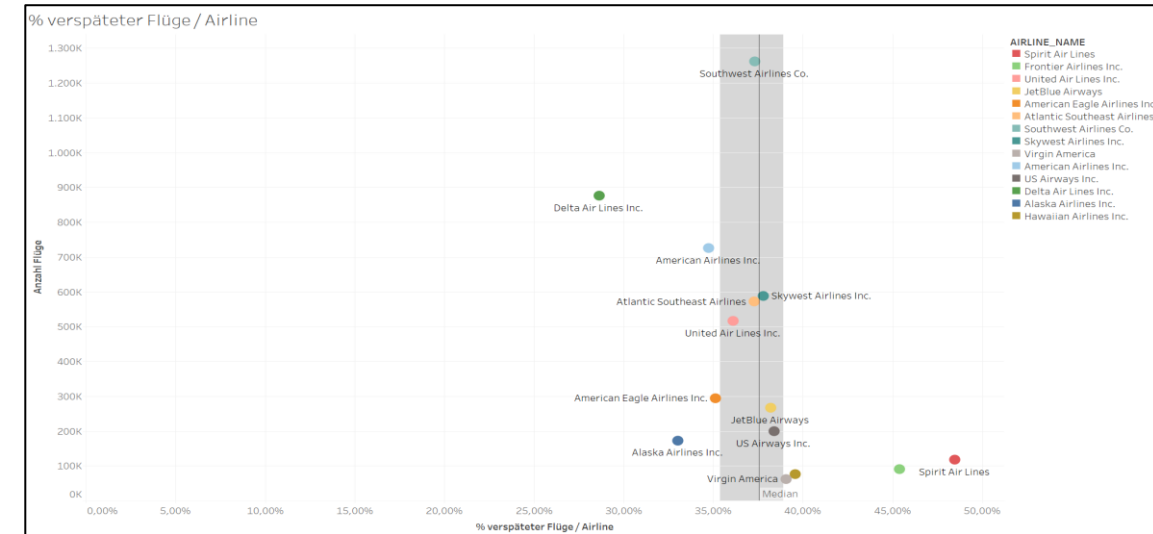
Analyse mit Tableau Desktop absolute Zahlen

- Grafiken sehr ähnlich
 - Southwest Airlines Co. immer führend
 - Grund Anzahl Flüge
- Ausnahme Annullierungen
 - American Eagle Airlines Inc. vergleichsweise viele Annullierungen
 - Delta Air Lines Inc. vergleichsweise wenige Annullierungen



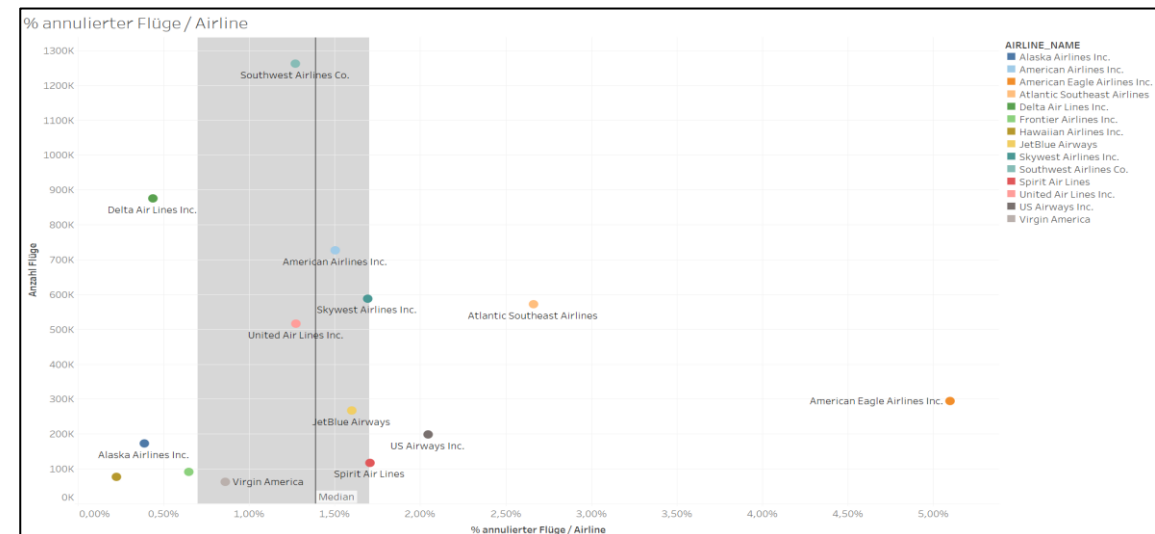
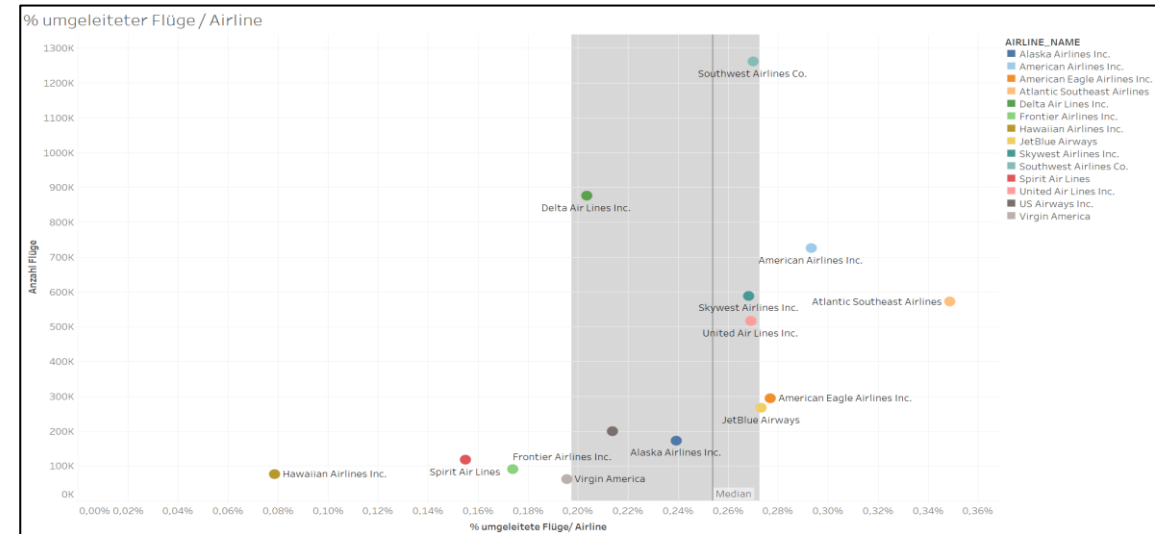
Analyse mit Tableau Desktop Prozentzahlen

- Hälfte der Airlines zwischen 35% und 38% der Flüge verspätet.
- Am meisten verspätete Airlines
 - Spirit Air Lines, 48% der Flüge
 - Frontier Airlines 45% der Flüge
- Am wenigsten verspätete Airline
 - Delta Air Lines Inc. 29% der Flüge
- Pünktlichste Airlines = verspätete Airlines mit umgekehrten Vorzeichen



Analyse mit Tableau Desktop Prozentzahlen

- Hälfte der Airlines zwischen 0.2% und 0.27% der Flüge umgeleitet
- Am wenigsten umgeleitete Airline
 - Hawaiian Airlines 0.08% der Flüge
- Am meisten umgeleitete Airline
 - Atlantic Southeast Airlines 0.35% der Flüge
- Hälfte der Airlines zwischen 0.7% und 1.7% der Flüge annulliert
- Am wenigsten annullierte Flüge
 - Hawaiian Airlines 0.22% der Flüge
- Am meisten annullierte Airline
 - American Eagle Airlines 5.1% der Flüge
- Delta Airlines Inc. 875'000 Flüge nur 0.44% annulliert



Zusammenfassung Analyse

- Southwest Airlines Co.
 - Bei absoluten Zahlen immer führend
 - Bei relativer Betrachtung immer in der Hälfte der Airlines
 - Grund Anzahl Flüge
- Spirit Air Lines
 - Am meisten verspätete/am wenigsten pünktliche Flüge
 - Wenig umgeleitete Flüge
 - Im Mittel bei annullierten Flügen
- American Eagle Airlines Inc
 - Relative viele annullierte Flüge
 - Ansonsten im Mittel
- Hawaiian Airlines
 - Wenigsten annullierten und umgeleiteten Flüge
 - Pünktlicher als die Hälfte der Airlines
 - Anzahl der Flüge sehr gering
- Delta Airlines Inc.
 - Am wenigsten verspätete/ am meisten pünktliche Flüge
 - Weniger Annullierungen als die Hälfte der Airlines
 - Im Mittel bei umgeleiteten Flügen
 - Anzahl der Flüge hoch

Fazit / Reflexion

- ETL-Prozess in eigenem Projekt angewendet
- Neue Tools kennengelernt
- Intensivierung von Python
- Anspruchsvoll und zeitintensiv mit unvorhergesehenen Problemen
 - VM-Abgleich, Datenmenge, unvollständige Datensätze, usw.
- Planung von Teamwork und individueller Arbeit
- Neues Wissen erlangt, steigende Lernkurve und **SPASS!**