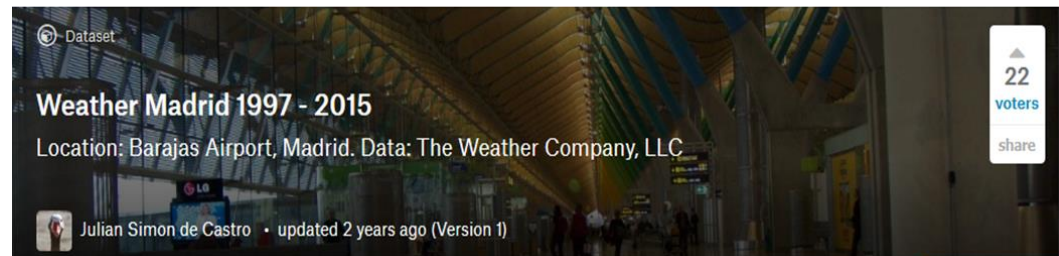


WETTER UND LUFTQUALITÄT VON MADRID

Inhalt

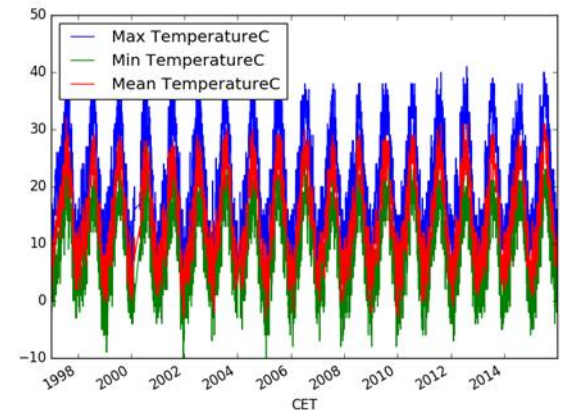
- Beschreibung der Daten
- Spark Befehlsabfolge
- Vergleich CSV und Parquet
- Verwendung von RDD
- Lessons Learned - Fazit



Beschreibung der Daten «Wetter»

Weather Madrid 1997 – 2015

- https://www.kaggle.com/juliansimon/weather_madrid_lemd_1997_2015.csv
- Files:
weather_madrid_LEMD_1997_2015.csv → 0.5 MB
- 6812 Zeilen
- pro Tag eine Angabe über
Temperatur, Feuchtigkeit, Luftdruck, Visibility,
Windgeschwindigkeit und -richtung, Niederschlag,
Bedeckung, Vorkommnisse, Wolkendecke



Beschreibung der Daten «Luftqualität»

Air Quality in Madrid (2001 - 2018)

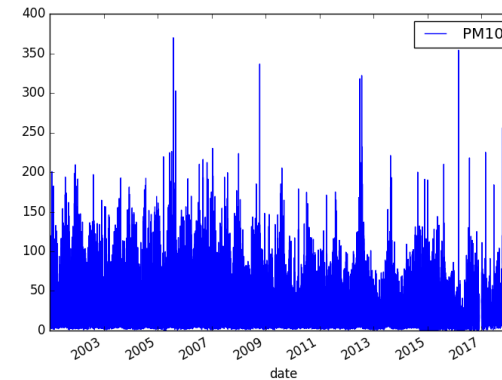
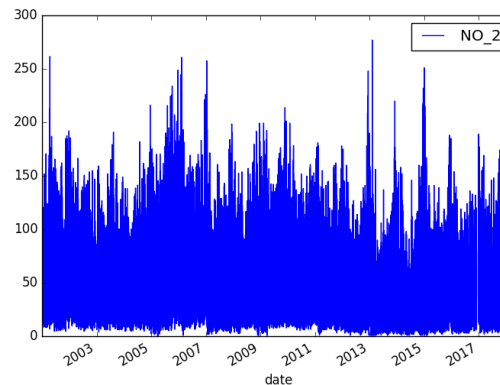
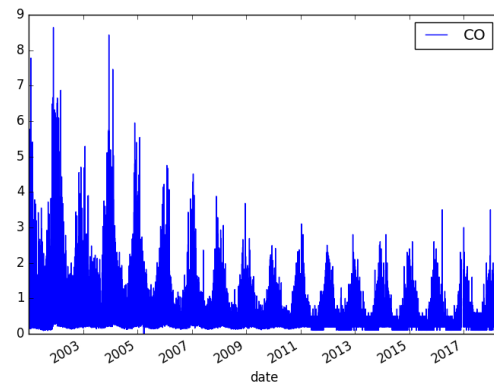
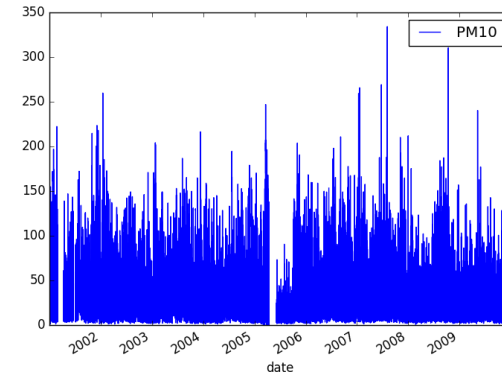
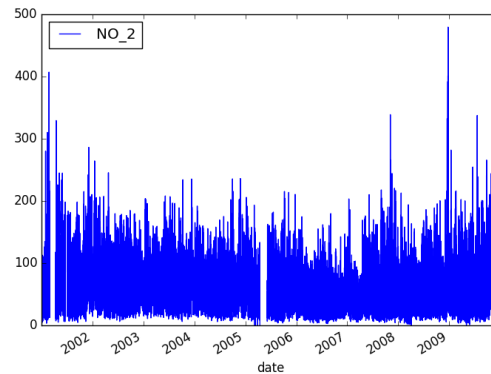
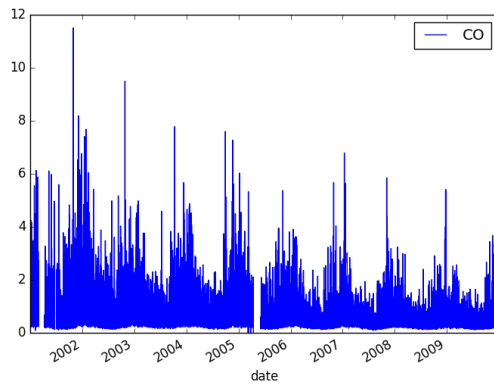
- <https://www.kaggle.com/decide-soluciones/air-quality-madrid>
- Files:
stations.csv – 24 Messstationen
madrid_2001.csv ... madrid_2018.csv → 500 MB
- 3.8 Millionen Tuples insgesamt
- pro Messstation jede Stunde eine Messreihe über Kohlenstoffmonoxid CO, Feinstaub PM10, Stickstoffdioxid NO2, Ozon O3, Nicht-Methan Kohlenwasserstoff NMHC, Benzene BEN und einige weitere

Luftqualität Kohlenstoffmonoxid / Stickstoffdioxid / Feinstaub zwei Stationen

Station:

28079007

← nur bis 2010 !



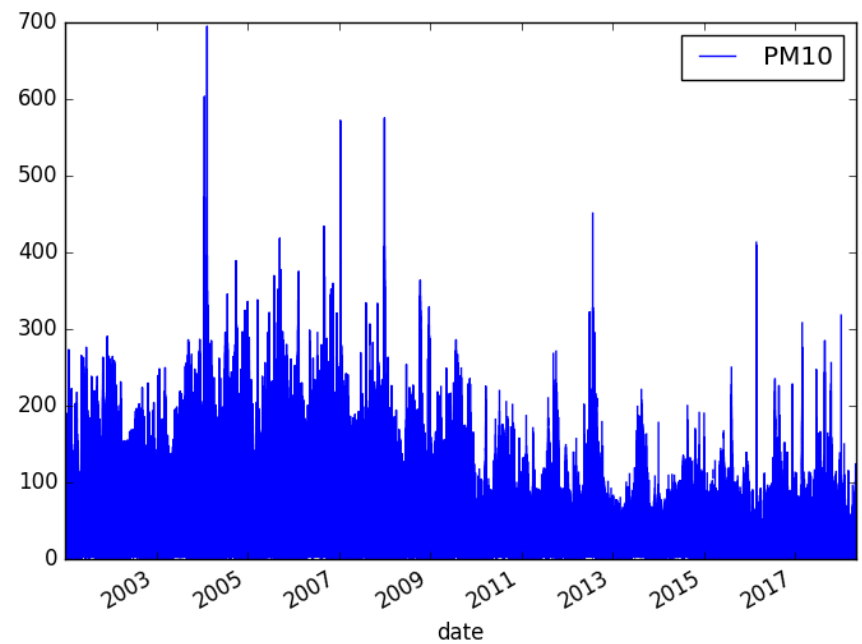
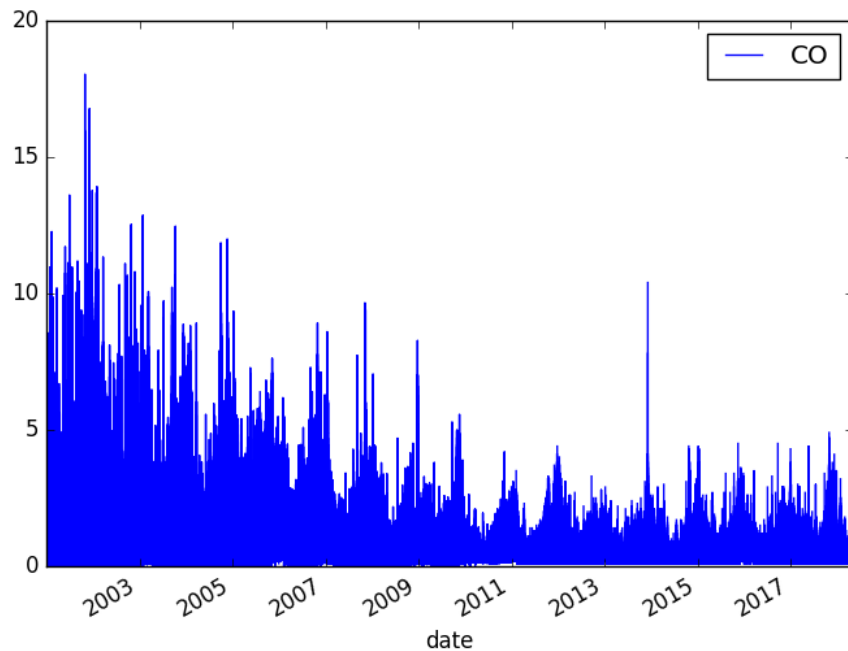
28079036

Vorgehen Spark - databricks

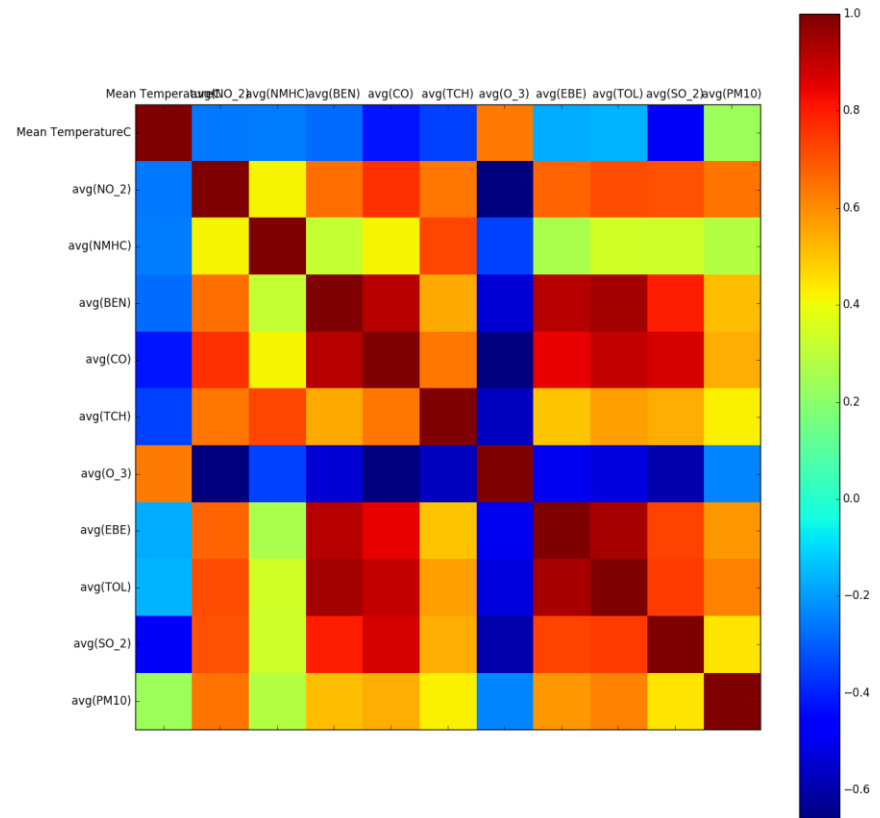
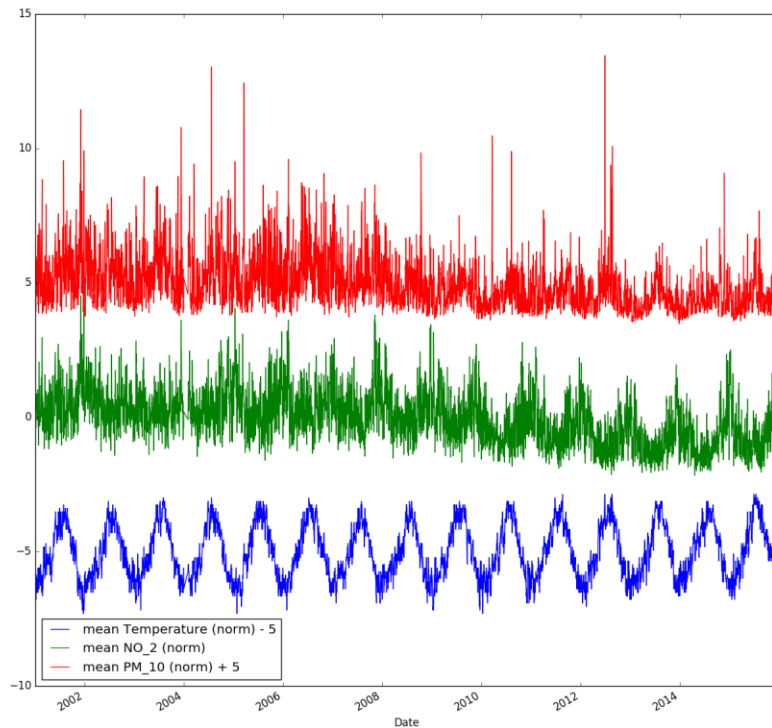
- Einlesen der CSV Daten:
- SQL-Abfrage:
- Daten ansehen:
- Grafisch Darstellung
- Daten zusammenführen:
- Reduktion auf gemeins. Attribute:
- Abspeichern in einem CSV-File:
- Abspeichern in Parquet-File:
- Format Zeitstempel anpassen:
- Grafiken mit gemittelten Daten:
- Join mit Wetterdaten:

```
sqlContext.read...  
sqlContext.sql('SELECT ... ')  
show()/count()/printSchema()  
toPandas()/plot(...)/display()  
unionByName() / unionAll()  
drop()  
df.write.csv(...)  
df.write.parquet(..)  
AVG()  
plot(...)  
join()
```

Kohlenstoffmonoxid und Feinstaub aller Stationen



Abhängigkeiten Schadstoffe – Temperatur?



Einsatz von RDDs

- Für die Reduktion der Datumsangaben in den Air Pollution Daten (3.8 Mio Datensätze) von Stunden auf Tage wurde der Dataframe auf ein RDD transformiert und über ein Mapping die Datumsangaben ersetzt.

```
1 # helper method to allow normalization of dates by "deleting" hours and minutes, i.e. reducing to day information
2 import datetime
3 from pyspark import sql
4
5 def replace_date (zeile):
6     datum = zeile["date"]
7     z = zeile.asDict()
8     z['date'] = datum.replace(hour=0, minute=0)
9     return sql.types.Row(**z)
10
```

```
1 airRDD = air_df.rdd
2 # normalize all timestamps to day only, i.e. overwrite hour and minute by 0
3 airRDD_new = airRDD.map(lambda x: replace_date(x))
4 # convert back to dataframe
5 air_df_new = spark.createDataFrame(airRDD_new)
6 air_df_new.createOrReplaceTempView("airRDD_new")
```


Vergleich CSV und Parquet

Join der Wetter Daten mit den Luftqualitätsdaten über das Datum.

Air

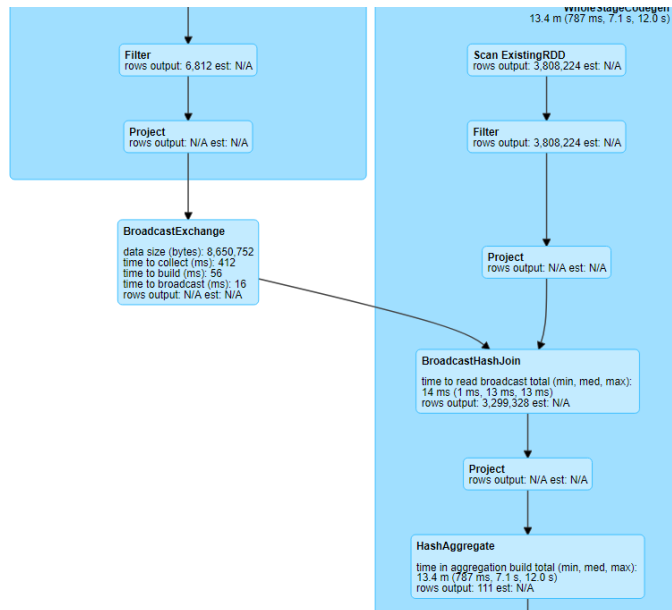
...			
datetime(2001, 8, 1, 0, 0)			
datetime(2001, 8, 1, 0, 0)			
...			
datetime(2001, 8, 1, 0, 0)			
...			

Weather (csv)

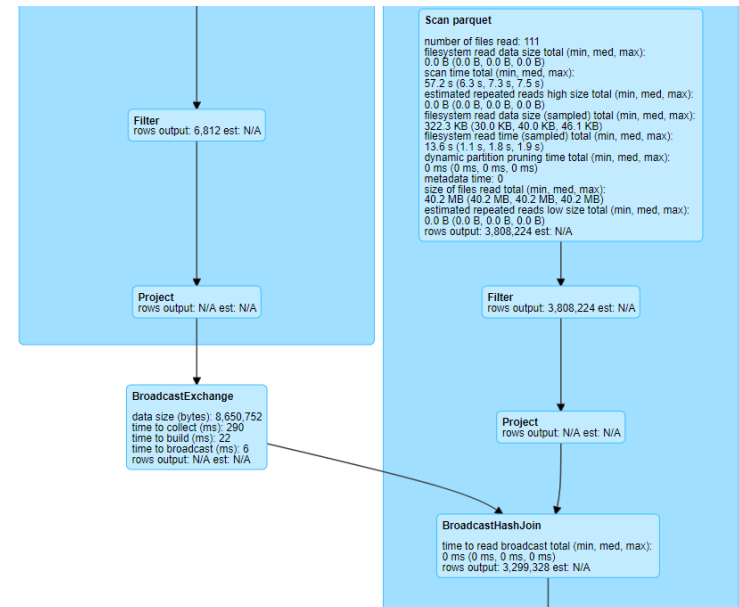
...			
datetime(2001, 8, 1, 0, 0)			
....			

24

Air als CSV



Air als Parquet



Lessons Learned - Fazit

- Einlesen der CSV Dateien (eine pro Jahr) wird sehr schön parallelisiert: 18 Dateien -> 36 Jobs
- Parquet mag keine Leerzeichen in Attributnamen. Umbenennen mit `df.withColumnRenamed(<old-name>, <new-name>).collect()`
- Darauf achten, wo man den Cache Befehl platziert, damit er zum richtigen Zeitpunkt Wirkung erzielt (Lazy Evaluation).
- Bei Caching andere Execution Pläne als ohne Caching
- Auch mit Caching (in memory) join mit Parquet immer noch rund Faktor 5 schneller (ohne Caching Faktor 15)
- Transformation von Spark Dataframes in Pandas Dataframes kann seeehr lange dauern. Folgender Trick kann helfen:
`spark.conf.set("spark.sql.execution.arrow.enabled", "true")`