

WHY IS DATA ENGINEERING SO HARD?

A LOOK AT THE CAUSE, SOLUTION AND FUTURE OF DATA ENGINEERING
DIFFICULTIES

ABOUT ME



Nearly 15 years
experience
dealing with data



Real-Time
Streaming



Predictive
Analytics



Sports Betting @
PhoenixHSL



Financial Services
@ Westpac /
Accenture

David Tout

- * Philosophy graduate
- * Strong opinions on data engineering
- * Wants you young kids off his lawn...



DATA ENGINEERING AS A PEOPLE PROBLEM

1. PROBLEM:
WE MAKE IT HARDER THAN IT NEEDS TO BE
2. SOLUTION:
WE NEED NAVIGATORS NOT HEROES
3. RESULT:
SELF-SERVICE DATA CONSUMPTION IS OUR GOAL





Engineers are
problem focused

Need to be more
people focused



Prefer our own ideas
(IKEA Effect)

Re-invent rather than
re-use and extend



Tools and technologies proliferate, littering
the landscape



DevOps and Data Science don't speak a
common language



Different goals & Different skillsets

WE ARE THE
PROBLEM!

CSV TO PARQUET CONVERSION METHODS

PANDAS & PYARROW

```
# csv_to_parquet.py

import pandas as pd
import pyarrow as pa
import pyarrow.parquet as pq

csv_file = '/path/to/my.tsv'
parquet_file = '/path/to/my.parquet'
chunksize = 100_000

csv_stream = pd.read_csv(csv_file, sep='\t', chunksize=chunksize, low_memory=False)

for i, chunk in enumerate(csv_stream):
    print("Chunk", i)
    if i == 0:
        # Guess the schema of the CSV file from the first chunk
        parquet_schema = pa.Table.from_pandas(df=chunk).schema
        # Open a Parquet file for writing
        parquet_writer = pq.ParquetWriter(parquet_file, parquet_schema, compression='snappy')
    # Write CSV chunk to the parquet file
    table = pa.Table.from_pandas(chunk, schema=parquet_schema)
    parquet_writer.write_table(table)

parquet_writer.close()
```

SPARK 2.0

```
1 from pyspark import SparkConf
2 from pyspark import SparkContext
3 from pyspark.sql import SQLContext
4
5 conf = SparkConf().setMaster("spark://bigdata-server:7077")
6 sc = SparkContext(conf=conf, appName="flightDataAnalysis")
7 sqlContext = SQLContext(sc)
8
9 #converts a line into tuple
10 def airlineTuple(line):
11     values = line.split(",")
12
13     return (
14         values[0], values[1], values[2], values[3], values[4], values[5],
15         values[10], values[11], values[12], values[13], values[14], values[15],
16         values[20], values[21], values[22], values[23], values[24], values[25])
17
18 #load the airline data and convert into an RDD of tuples
19 lines = sc.textFile("hdfs://localhost:9000/user/bigdata/airline/input/airline.csv")
20
21 #convert the rdd into a dataframe
22 df = sqlContext.createDataFrame(lines, ['Year', 'Month', 'DayofMonth', 'DayOfWeek',
23                                         'CRSArrTime', 'UniqueCarrier', 'FlightNum',
24                                         'CRSElapsedTime', 'AirTime', 'TailNum',
25                                         'Distance', 'TaxiIn', 'TaxiOut', 'WheelTimeIn',
26                                         'WheelTimeOut', 'CarrierDelay', 'WeatherDelay',
27                                         'LateAircraftDelay'])
28
29 #save the dataframe as a parquet file in HDFS
30 df.write.parquet("hdfs://localhost:9000/user/bigdata/airline/input/airline.parquet")
```


THE NEED FOR NAVIGATORS NOT HEROES



THE WORLD *WILD* WEB



Web design used to be difficult



Fancy tools like wix, or squarespace make many web designer roles redundant.



Web roles still exist, but for much more involved sites



48% chance of automation by AI for Computer Programmers, according to a 2013 study



However, **Database Administrators** scored only 3%

Data Engineers are probably at **3-8%** risk



These figures will grow, year on year, and eventually accelerate.



If your job is the 2019 equivalent of Web Design, you stand a good chance to lose your job in the next two decades.

AI IS
COMING
FOR
YOU_(R JOB)

SELF SERVICE PLATFORMS & ENVIRONMENTS



Let others work...



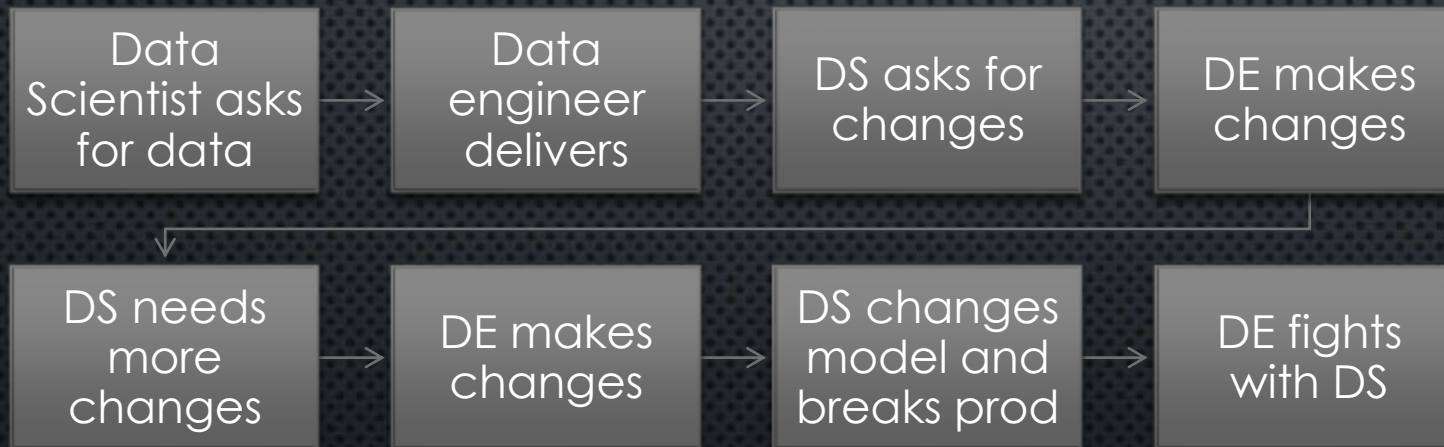
While you be yourself...

PREDICTING WEATHER PATTERNS ON VENUS

- A HYPOTHETICAL SCENARIO FOR THE INTERACTION OF A DATA SCIENTIST WITH TWO TYPES OF DATA ENGINEERS (NAVIGATOR VS HERO)



HERO DATA ENGINEER APPROACH





NAVIGATOR DATA ENGINEER APPROACH



WHAT HAVE WE LEARNED?

- DESIGN SELF-SERVICE
PLATFORMS
- SCALE BEYOND YOURSELF
- MAKE YOURSELF
REDUNDANT
- TEACH & LEARN
- ENABLE OTHERS





ARE YOU A
NAVIGATOR
OR A HERO?