

Data Engineering mit Apache Kafka

Projektbericht
von

Johannes Weber

Matrikelnummer: 11010021

und

Julian Ruppel

Matrikelnummer: 11010020

09.02.2018

SRH Heidelberg
Fakultät für Information, Medien und Design
Big Data und Business Analytics

Dozent
Frank Schulz

Inhaltsverzeichnis

1	Einleitung	1
1.1	Aufgabe und Ziel	1
2	Tools	3
3	Lösungsansatz	4
3.1	Architekturentscheidungen	4
3.2	Data Ingestion	4
3.2.1	Data Ingestion via CSV Datei	4
3.2.2	Data Ingestion via OData Schnittstelle	4
3.3	Data Storage	4
3.4	Data Retrieval	5
3.4.1	Data Retrieval mit Apache Zeppelin	5
3.4.2	Data Retrieval mit Jupyter	5
4	Fazit	6
	Abbildungsverzeichnis	ii
	Abkürzungsverzeichnis	iii

Kapitel 1

Einleitung

1.1 Aufgabe und Ziel

Im Rahmen dieses Projektes war es die Zielstellung sich mit Data Ingestion, Data Storage sowie Data Retrieval vertraut zu machen.

Data Ingestion ist die Beschaffung der Daten. Dies kann entweder mit Hilfe eines Data Streams erfolgen oder einer statischen Datenquelle - also eine einfache Datei die lokal auf einem Rechner angelegt wird und Daten beinhaltet wie z. B. eine Comma-separated values (CSV) oder JavaScript Object Notation (JSON) Datei. Unter einem Data Stream versteht man einen kontinuierlichen Datenstrom wie z. B. die Erstellung von immer wieder neuen Twitter Nachrichten. Ein wichtiges Merkmal eines Data Streams ist, dass man nicht vorhersehen kann wann der Datenstrom zu Ende ist - er könnte theoretisch unendlich sein. Im Falle von einem Datenstrom von Twitter Nachrichten ist es nicht abzusehen wann jemals die letzte Twitter Nachricht geschrieben wird. Für unsere Aufgabe ist darauf zu achten, dass der Datenstrom über eine API öffentlich zugänglich ist und immer auf dem aktuellsten Stand gehalten wird.

Data Storage ist die Speicherung der Daten. Hierbei wurde uns lediglich die Anforderung gestellt, dass wir für die Speicherung die Streaming Plattform Apache Kafka verwenden. Des Weiteren war es uns gestattet die Daten in einer relationalen Datenbank, NoSQL Datenbank oder mit Spark Streaming speichern, falls nur die Nutzung von Apache Kafka unsere Anforderungen nicht genügt.

Data Retrieval ist die Beschaffung der Daten aus einer Datenbank mit SQL Abfragen und die abschließende Ausgabe der Ergebnisse in Form von Tabellen oder einfachen Visualisierungen. Die Visualisierung der Daten sollte in einem virtuellen Notebook erfolgen. Als virtuelles Notebook durften wir uns entscheiden zwischen Apache Zeppelin oder Jupyter.

Unsere Aufgabe ist es ein geeignetes Szenario für die Bewältigung dieser Aufgabe zu finden und umzusetzen.

Wir entschieden uns für unser Szenario die Socrata Application Programming Interface (API) von NYC Open Data zu nutzen. NYC Open Data ermöglicht es allen "New Yorker" und somit auch der ganzen Welt sogenannte Open Data also frei zugängliche Daten einfach zu konsumieren.¹.

NYC Open Data ermöglicht es uns sowohl einen kontinuierlichen Data Stream als auch eine statische CSV Datei zu konsumieren. Dank diesem Umstand entschieden wir uns im Rahmen dieses Projektes beide Möglichkeiten umzusetzen und zu vergleichen. Die Data Ingestion mit der CSV Datei setzte Julian Ruppel mit der Programmiersprache Java um und den kontinuierlichen Datenstrom über die Socrata API wurde von Johannes Weber mit der Programmiersprache Python ausgelesen.

Auch bei der Data Retrieval entschieden wir uns dafür sowohl Apache Zeppelin als auch Jupyter zu nutzen und zu vergleichen. Die Beschaffung und Auswertung der Daten mit Apache Zeppelin übernahm Julian Ruppel. Johannes Weber bereitete die Daten mit der Programmiersprache Python auf und visualisierte sie in einem Jupyter Notebook.

Kapitel 2 - *Tools* beschäftigt sich detaillierter mit den verwendeten Tools, Frameworks und Programmiersprachen und Kapitel 3 - *Lösungsansatz* erläutert das gewählte Szenario mit den Unterkapitel Data Ingestion, Data Storage und Data Retrieval sowie der verwendeten Architektur.²

¹?

²Einleitung von Johannes Weber

Kapitel 2

Tools

Kurzworstellung aller Tools, Frameworks und Sprache

Kapitel 3

Lösungsansatz

3.1 Architekturentscheidungen

Lambda vs. Kappa

Warum Lambda? stream ist unendlich. Ergebnisse/Aggregationen zwischenspeichern Aggregationen mit den Live Daten abmischen! Nachteile googlen Quelle verweisen

Warum Kappa? Daten sind endlich aber ausreichend für unsere Bedürfnisse. Haben eh nur Pseudostream. Vorteile googeln

<https://www.confluent.io/blog/simplest-useful-kafka-connect-data-pipeline-world-thereabouts-part-1/> Grafik einbetten

3.2 Data Ingestion

3.2.1 Data Ingestion via CSV Datei

Da muss was stehen

3.2.2 Data Ingestion via OData Schnittstelle

Da muss was stehen

3.3 Data Storage

Data Storage

3.4 Data Retrieval

3.4.1 Data Retrieval mit Apache Zeppelin

Da muss was stehen

3.4.2 Data Retrieval mit Jupyter

Da muss was stehen

Kapitel 4

Fazit

Abbildungsverzeichnis

Abkürzungsverzeichnis

JSON JavaScript Object Notation

CSV Comma-separated values

API Application Programming Interface