

Studienarbeit - Entwurf

Ausarbeitung InfluxDB und Wetterdaten

Erstellt von: Henrik Mertens Hatzfelder Str 25 33104 Paderborn

Prüfer: Prof. Dr. Ulrich. Reus

> Eingereicht am: 19. Juni 2022

Inhaltsverzeichnis

A	Abkürzungsverzeichnis I				
A	bbild	ngsverzeichnis	\mathbf{V}		
Ta	abell	verzeichnis	VI		
Li	sting	erzeichnis V	'II		
1 Einleitung					
	1.1	ielsetzung	1		
	1.2	ufbau und Vorgehensweise	1		
2	Gru	dlagen	2		
	2.1	Time Series Data	2		
	2.2	ufbau von Time Series Datenbanken	3		
	2.3	Unterschiede zwischen Time Series und relationalen Datenbanken $\ . \ .$	3		
	2.4	Verbreitete DBMS	4		
	2.5	Intwicklungsumgebung	4		
		.5.1 Docker und Docker-Compose	5		
		.5.2 Python und Jupyter Notebooks	5		
3	Infl	:DB	7		
	3.1	nfluxDB Installation	7		
	3.2	Daten einfügen	8		
		.2.1 Line Protokoll	8		
	3.3	Daten abrufen	9		
	3.4	Ython library	9		
	3.5	Oaten verarbeiten und visualisieren	10		
	3.6	reitere InfluxDB-Funktionen	10		
	3.7	Vetterdaten	10		
		.7.1 Wetterdaten Aufbau	10		
		.7.2 Wetterdaten abrufen	11		
		.7.3 Wetterdaten Verarbeiten	12		
4	\mathbf{Zus}	nmenfassung	14		

Anhang	15
Quellenverzeichnis	23
Ehrenwörtliche Erklärung	25

Abkürzungsverzeichnis

API	Application Programming Interface.
CLI	Command Line Interface.
CRUD	Create, Read, Update, Delete.
CSV	Comma-separated values.
DWD	Deutscher Wetterdienst.
HTML	Hypertext Markup Language.
HTTP	HyperText Transfer Protocol.
IOT	Internet of Things.
RDBMS	Relational Database Management System.
TSDB	Time Series Database.
UI	User Interface.

Abbildungsverzeichnis

Abbildung 1:	DB-Engines Ranking	17
Abbildung 2:	InfluxDB Dashboard	18
Abbildung 3:	InfluxDB Load Data Source	19
Abbildung 4:	InfluxDB Load Data Bucket	20
Abbildung 5:	InfluxDB Load Data Bucket hinzufügen	20
Abbildung 6:	InfluxDB Load Data API Tokens	21
Abbildung 7:	InfluxDB Load Data API Token hinzufügen	22

Tabellenverzeichnis

Tabelle 1:	Bedeutung der CSV Felder	12
Tabelle 2:	Bedeutung Stationstabellen Felder	13

Listingverzeichnis

Listing 1:	InfluxDB Line Protokoll Quelle: Influxdata (2022b)	8
Listing 2:	Wetterdaten CSV Quelle: DWD Wetterdaten CSV Download $~$.	11
Listing 3:	Influx DB Python Daten schreiben Quelle: Influx data (2022d) $$. $$.	16
Listing 4:	Influx DB Python Daten lesen Quelle: Influx data (2022d) $\ \ . \ .$	16

1 Einleitung

1.1 Zielsetzung

Das Ziel dieser Arbeit ist es eine Einführung in die Funktion von Time Series Database (TSDB) zu geben. Dazu wird beispielhaft an den Temperaturdaten des Deutschen Wetterdienstes (DWDs) gezeigt wie einer TSDB gearbeitet werden kann. Als Datenbank wird InfluxDB genutzt.

1.2 Aufbau und Vorgehensweise

Im Grundlagenteil dieser Arbeit werden die Grundlagen von Time Series Data und TSDB erläutert. Des Weiteren wird erklärt was eine TSDB ist und wie diese sich von Relational Database Management System (RDBMS) unterscheidet. Außerdem wird kurz die für die Entwicklungsumgebung notwendige Software Docker und jupyter Notebooks erklärt.

Im nächsten Kapitel wird gezeigt wie mit InfluxDB gearbeitet werden kann. Dazu wird das für InfluxDB wichtige Line Protokoll vorgestellt. Außerdem wird gezeigt wie die abgerufenen Daten mit InfluxDB visualisiert werden können. Als Daten werden die Wetterdaten das DWD genutzt welche auch in diesem Kapitel erklärt und vorstellt werden.

Im letzten Artikel folgt eine kurze Zusammenfassung der Arbeit und der wichtigsten Informationen zu TSDB

2 Grundlagen

In diesem Kapitel werden die Grundlagen von TSDB und Time Series Data erklärt. TSDB gehören zu den NoSQL Datenbanken und sind besonders darauf optimiert, mit Time Series Data zu arbeiten. Dadurch können die große Mengen an Time Series Data verarbeiten, durchsuchen und speichern.¹

2.1 Time Series Data

Um TSDB zu verstehen, muss als erstes geklärt werden, was Time Series Data überhaupt ist und wie sie sich von anderen Daten unterscheiden. Wie der Name schon impliziert, ist Time Series Data eine Reihe von Daten, die über einen Zeitraum gesammelt worden sind. Es wird also nicht nur der Endwert aufgezeichnet, sonder die Veränderung über einen Zeitraum. Diese Daten können z.B. Servermetriken, Netzwerkdaten, Internet of Things (IOT) Sensordaten, Ereignisse, Klicks, Marktgeschäfte und viele andere Arten von Daten sein. Time Series Data können gut daran erkannt werden, dass die Zeit eine wichtige Achse bei der Darstellung der Werte ist.²

Manchmal ist es nicht notwendig, alle Daten zu erfassen. Zum Beispiel wird in vielen Anwendungen nur der letze Login gespeichert. Mehr ist auch für die Funktion nicht notwendig. Allerdings können zusätzliche Informationen gewonnen werden, wenn nicht nur der letze Datenpunkt sondern die Veränderung über einen Zeitraum aufgezeichnet wird. So kann zum Beispiel festgestellt werden, wie oft und wann sich der Kunde einloggt und ob es dabei ein Muster gibt. Anhand dieser Daten können Kunden dann kategorisiert werden.³

Eine weitere Besonderheit von Time Series Data ist, dass sie sich nicht verändert. Wenn die Daten einmal erfasst wurden wird an ihnen nichts mehr verändert. Es werden nur neue Daten hinzugefügt.⁴

¹vgl. ComputerWeekly.de, Redaktion (2021)

 $^{^{2}}$ vgl. Dix, Paul (2021), S. 1 ff.

 $^{^3\}mathrm{vgl.}$ Data-Science-Team (2020)

 $^{^{4}}$ vgl. Fangman, Sam (2019)

2.2 Aufbau von Time Series Datenbanken

TSDB sind darauf ausgelegt Key Value paare zu Speichern. Der Key in einem Key Value Datensatz ist eine Wert über den die Value referenziert wird. Im Value Teil werden die Daten zum dazugehörigem Key gespeichert. Der Wert der Value kann ein primitiver Datentyp sein oder auch ein Objekt das in einen primitiven Datentyp umgewandelt worden ist. 5

Ein Datenpunkt in einer Time Series Database besteht aus mehreren key Value Paaren. Einige Dieser Key Value Paare sind sogenannte Tags. Diese Tags sind Werte die sich zwischen den Datenpunkten nicht ändern wie zum Beispiel die Position eines Sensors oder die Kundenummer eines Kunden der sich gerade eingeloggt hat. Anhand dieser Tags können die Datenpunkte durchsucht werden. Die eigentlichen Messwerte des Sensors oder andere Daten die erfasst werden sollen werden auch als Key Value Paar gespeichert. Zum Beispiel wird als Name Temperatur und als Wert 25,2 angeben. Ein Datenpunkt kann mehrere Messwerte haben. Außerdem wird jeder Datenpunkt mit einem Timestamp versehen nach welchem er Indexiert wird.⁶ Eine gute Veranschaulichung wie die Daten in einer TSDB aufgebaut sind zeigt das InfluxDB Line Protokoll in Listing 1.

2.3 Unterschiede zwischen Time Series und relationalen Datenbanken

Um Time Series Data zu speichern, ist es nicht unbedingt erforderlich, eine TSDB zu nutzen. Auch relationale Datenbanken können Time Series Data speichern. Einer der wichtigsten Unterschiede zwischen einer TSDB im Gegensatz zu einem RDBMS ist es, dass kein Datenbank Schema benötigt wird. Wenn Time Series Daten in einer Rationalen Datenbank geschrieben werden sollen, müssen erst entsprechende Tabellen angelegt werden, in denen die Daten immer im gleichen Format abgelegt werden müssen. Im Gegensatz dazu können in einer TSDB die Daten einfach schemafrei in die Datenbank geschrieben werden. Ein weiterer Vorteil ist, dass TSDB im Gegensatz zu relationalen Datenbanken besser und einfacher Skaliert werden können.⁷

 $^{^{5}}$ vgl. Seeger, Marc (2009)

 $^{^6 \}mathrm{vgl.}$ hazelcast (2022)

 $^{^{7}}$ vgl. Influxdata (2021)

Aber TSDB haben nicht nur Vorteile. Wie in Abb. 1 zu sehen, sind sie viel weniger verbreitet als nicht zeit basierte Datenbank Systeme. Dadurch gibt es viel weniger Entwickler die Erfahrungen mit TSDB haben. Auch auch das Ökosystem um die Datenbank ist deutlich kleiner. Außerdem sind RDBMS dadurch, dass es sie viel länger gibt, sehr stabil und sehr gut unterstützt.⁸

RDBMS arbeiten nach dem Create, Read, Update, Delete (CRUD) Prinzip, welches für Time Series Data nicht optimal ist. Auf Time Series Data werden keine Update Befehle durchgeführt, da neue Daten immer nur angehängt werden. Auch das Löschen von Daten wird nicht sehr häufig durchgeführt und im Gegensatz zu RDBMS meistens gleichzeitig auf einer großen Menge an Datensätzen. TSDB sind auf diese Besonderheiten optimiert und daher besser dafür geeignet, mit Time Series Data zu arbeiten und weisen auch eine höhere Performance auf.⁹

2.4 Verbreitete DBMS

Aktuell gibt es wie in Abb. 1 zu sehen, einige beliebte Mulit-Model Datenbanken, die als TSDB genutzt werden können. So können die Datenbanken wie MongoDB, Redis, Teradata und Couchbase mit Time Series Daten arbeiten. Die erste reine TSDB im Ranking ist InfluxDB auf Platz 29.¹⁰

Allerdings haben Datenbanken, die nur auf das verarbeiten von Time Series Data ausgelegt sind, deutliche Performance Vorteile gegenüber Multi Model Datenbanken. In einem Vergleich von InfluxDB und MongoDB, hat InfluxDB eine 2,4 mal bessere Schreibperformance als MongoDB und ist beim Lesen sogar 5,7 mal schneller. InfluxDB benötigt außerdem 20 mal weniger Speicherplatz auf der Festplatte, um die gleiche Menge an Daten zu speichern.¹¹

2.5 Entwicklungsumgebung

Um mit InfluxDB zu arbeiten wird eine Umgebung zum ausführen von Docker Containern benötigen, in welchen wir InfluxDB und Jupyter Notebooks betreiben. Der

 $^{^{8}}$ vgl. Influxdata (2021)

⁹vgl. Influxdata (2021)

 $^{^{10}}$ vgl. solid-IT-gmbh (2022)

 $^{^{11}\}mathrm{vgl.}$ Hajek Vlasta Pour Ales, Kudibal Ivan (2019)

eigentliche Code wird dann in Jupyter Notebooks mit Python entwickelt. Die Grundlagen über die Eingesetzen Tool und Techniken werden grob in diesem Kapitel erläutert.

2.5.1 Docker und Docker-Compose

Docker ist eine Software für das erstellen und verwalten von Containern. Mit Docker ist es möglich Anwendungen samt ihrer Umgebung in einer Einheit zusammenzufassen, so das diese einfach auf anderen System ausgeführt werden können. Dabei hat jeder Container ein eigenes Dateisystem und ein eigens Betriebssystem. Allerdings teilen sich Container und Hostsystem den Kernel des Hostsystems. Dadurch hat diese Art der Virtualisierung deutlich weniger Overhead als andere Virtualisierungstechniken. Zusätzlich wird das Betriebssystem innerhalb des Containers maximal reduziert so das nur noch benötigte Komponenten vorhanden sind. Wichtig ist es das immer nur möglichst eine Anwendung in einem Container zu finden ist. Durch die Virtualisierung sind die einzelnen Container voneinander getrennt.¹²

Allerdings bestehen einige Anwendungen aus mehreren Komponenten, diese können durch mehrere Docker Container abgebildet werden. Um die Verwaltung von mehreren Container zu erleichtern kann Docker-Compose genutzt werden. Mithilfe von Docker Compose können größere Umgebungen in einem Compose File verwaltete werden. Hier werden die Umgebungsvariablen, Container Image oder Dockerfiles, Ports, Storage und weiteres in einer Datei definiert. Mithilfe dieser definition kann Docker Compose eine komplexe Umgebung mit nur einem Befehl initialisieren.¹³

2.5.2 Python und Jupyter Notebooks

Python ist eine universelle Prozedurale und Imperative Programmiersprache die 1994 in der ersten Version veröffentlicht wurde. Der Name ist eine Huldigung an Monty Python und wurde nicht nach einer Schlange benannt, auch wenn das Logo eine Schlange ist. Python ist unter der freien PFS Lizenz lizenziert wodurch es auch in kommerziellen Anwendung genutzt werden kann. Python ist eine Interpretierte Sprache. Das heißt das sie nicht zu einer ausführbaren Datei kompiliert wird sondern von einem Interpreter

 $^{^{12}}$ vgl. Stender, Daniel (2020), S. 54 ff.

 $^{^{13}\}mathrm{vgl.}$ Stender, Daniel (2020), S. 151 ff.

interpretiert wird. Außerdem ist Python eine unter Programmieranfängern sehr beliebte Sprache die auch sehr viel in den Bereichen DataSience, DeepLerning, Naturwissenschaften und Linux Systemprogrammierung eingesetzt wird.¹⁴

Jupyter Notebooks ist eine Webbasierte Open-Source Anwendung mit dem Ziel Code in den Sprachen Python, R, und Julia einfach zu schreiben, bearbeiten, auszuführen und einfach zu teilen. Ein Notebook besteht immer aus Zellen. Eine Zelle kann Code oder Markdown Formatierten Text anzeigen. Jede Zelle kann einzeln ausgeführt werden. Dadurch kann ein Programm sehr einfach und verständlich dargestellt und erklärt werden. Es ist auch möglich neue Notebooks und Dateien im Webinterface von Jupyter Notebooks selbst anzulegen.¹⁵

¹⁴vgl. Stender, Daniel (2020), S. 66 ff.

¹⁵vgl. Silaparasetty, Nikita (2020), S. 91 ff.

3 InfluxDB

InfluxDB ist eine in Go geschriebene open source TSDB, die darauf ausgelegt ist, mit einer großen Menge an Time Series Data zu arbeiten.¹⁶ Im weiterem verlauf dieses Kapitels wird am Beispiel von Wetterdaten gezeigt, wie mit InfluxDB gearbeitet wird. InfluxDB stellt für die Integration in eigene Anwendungen ein HyperText Transfer Protocol (HTTP) Application Programming Interface (API) zur Verfügung, für die es in vielen Programmiersprachen Client Librarys gibt. Außerdem wird ein Webinterface und ein Command Line Interface (CLI) bereitgestellt.¹⁷

3.1 InfluxDB Installation

Bevor InfluxDB genutzt werden kann muss es als erstes installiert werden. Am einfachsten ist dies über Docker möglich. Dazu ist es notwendig das Docker und Docker Compose auf dem System installiert sind. Mit Docker Desktop lassen sich die beide Tools am einfachsten installieren. Im Anhang dieser Arbeit befindet sich im Ordner Docker eine Docker Compose Datei mit dem Namen docker-compose.yml. Zum Starten der benötigten Container ist es am einfachsten mit einem Terminal (Powershell, xterm usw.) in den Docker Ordner zu wechseln und den Befehl docker compose up -d auszuführen. Jetzt beginnt docker damit die notwendigen Images herunterzuladen und zu bauen. Wenn der Befehl ohne Fehler ausgeführt worden ist InfluxDB installiert worden und kann über die URL: http://localhost:8086 aufgerufen werden. Die Login Daten sind als Umgebungsvariable in Docker Compose definiert und lauten: admin e1LjSYaFbzbJeIBC.

Außerdem wurde mit diesem Befehl auch ein Jupyter Notebook in einem Docker Container gestartet. Auf diesen Container kann über die URL: http://localhost:8888/ zugegriffen werden. Das Passwort lautet fhdw. Im Ordner DWD befinden sich die Notebooks mit dem in dieser Arbeit beschrieben Code.

 $^{^{16}}$ vgl. solid-IT-gmbh (2022)

 $^{^{17}\}mathrm{vgl.}$ Influx
data (2022a)

3.2 Daten einfügen

In InfluxDB werden Daten immer in Buckets gespeichert. Um Daten hochzuladen, muss zunächst ein Bucket angelegt werden. Dazu gibt es zwei Möglichkeiten. Die Einfachste ist es, über das Web User Interface (UI) von InfluxDB einen neuen Bucket anzulegen. Dazu muss nach dem Login der Navigationspunkt Data und der Reiter Buckets ausgewählt werden. Hier kann dann mit dem Button Create Bucket ein neuer Bucket angelegt werden. Bei dem Anlegen kann noch eine Lebensdauer für die Daten ausgewählt werden, nach welcher die jeweiligen Datenpunkte gelöscht werden.¹⁸

3.2.1 Line Protokoll

Daten werden immer nach dem InfluxDB Line Protokoll formatiert an die Datenbank gesendet. Das Protokoll ist wie in Listing 1 dargestellt aufgebaut. Im ersten Teil des Line Protokolls wird der Name der Messreihe angegeben. Das kann zum Beispiel der Name des Sensors sein oder der Ort an dem der Messwert genommen wurde. Wichtig ist, dass Groß und Kleinschreibung beachtet werden muss und Unterstriche nicht genutzt werden dürfen. Sonderzeichen müssen mit einem \ maskiert werden. Nach dem Namen kommen getrennt durch ein Komma die Tags der Messung. Tags werden indexiert und dazu genutzt, um Messwerte zu durchsuchen. Tags werden als Key Value Paar angegeben. Hier sollen Metadaten wie zum Beispiel der Standort des Sensors oder der Name des Servers eingetragen werden, zu dem die Datenpunkt/e gehören. Die eigentlichen Werte sind mit einem Leerzeichen von den Tags abgegrenzt und bestehen aus durch Kommas getrennten Key Value Feldern. Der letzte Wert einer Zeile ist der Unix Timestamp in Millisekunden. In einer Datei oder Anfrage kann es mehrere Zeilen mit Daten geben.¹⁹

isting 1: InfluxDB Line Protokoll Quelle: Influxdata (2022b)							
measurementName	,tagKey=tagValue	fieldKey="fieldValue"	1465839830100400200				
Measurement	Tag set	Field set	Timestamp				

Die im Line Protokoll formatierten Daten können jetzt entweder mithilfe eines Rest Requests oder des InfluxDB CLI in die Datenbank übertragen werden. Um diese An-

 $^{^{18}\}mathrm{vgl.}$ Abb. 2, Abb. 3, Abb. 4, Abb. 5

¹⁹vgl. Influxdata (2022b)

fragen zu autorisieren muss ein API Token mitgesendet werden.²⁰ Um einen Token zu bekommen, kann dieser entweder über das Webinterface, die CLI oder über die API angelegt werden. Der einfachste Weg ist es, den Token über das Webinterface anzulegen. Dazu wird wie beim Anlegen eines Buckets zunächst der Menüpunkt Data ausgewählt und anschließend der Reiter API Tokens. Mit einem Klick auf Generate API Token kann dann ein API Token erstellt werden.²¹ Dabei kann zwischen einem All-Access token und einem Read/Write token ausgewählt werden. Mit dem All Access Token kann auf alles zugegriffen werden. Mit einem Read/Write Token kann wie in Abb. 7 zu sehen, ausgewählt werden, auf welchen Bucket geschrieben oder gelesen werden kann.²²

3.3 Daten abrufen

3.4 Python library

Um nicht selbst eigene API Anfragen über die Rest API schreiben zu müssen gibt es für Python und andere Sprachen fertige Bibliotheken.²³ Bevor mit der Client Library gearbeitet werden kann muss diese als erste Installiert und importiert werden. Am besten wird zur Installation der Python Paketmanager pip genutzt. Mit dem Befehl pip install influxdb-client werden alle benötigten Packte installiert.

In Listing 3 kann an einem Beispiel gesehen werden wie Daten mithilfe von Python in die Datenbank geschrieben werden. Ein ähnliches Beispiel findet sich ausführbar und detailliert beschreiben im Jupyter Notebook. In der ersten beiden Zeile des Codes wird der InfluxDB Client importiert. Danach werden in Zeile vier bis sieben die benötigten Daten bucket, Organisation, Token und URL in Variablen geschrieben. Was in die Variablen eingetragen werden muss wird im Kapitel Daten einfügen beschreiben.

In Zeile neun bis 13 wird der Client mit den hinterlegen Variablen initialisiert. Danach folgt in Zeile 15 die Initialisierung des write clients. Als nächstes muss ein Datenbank erstellt werden der in die Datenbank geschrieben werden kann. Dazu wird in Zeile 17 bis

²⁰vgl. Influxdata (2022f)

²¹vgl. Abb. 2, Abb. 3, Abb. 6, Abb. 7

 $^{^{22}}$ vgl. Influxdata (2022c)

 $^{^{23}}$ vgl. Influxdata (2022e)

19 der Datenpunkt erstellt. An diesen Datenpunkt kann eine beliebige Menge an Tags und Datenfeldern über die Methoden tag und field angehängt werden. Hier im Beispiel wird der Datenpunkt my_measurement mit dem Tag location angelegt, welcher den Wert Paderborn hat, und dem Datenpunkt temperature mit dem Wert 25,3 angelegt. In der letzten Zeile wird dann der write client dazu genutzt um diesen Datenbank an die Datenbank zu senden.

3.5 Daten verarbeiten und visualisieren

3.6 weitere InfluxDB-Funktionen

3.7 Wetterdaten

3.7.1 Wetterdaten Aufbau

Die Wetterdaten des DWD können über den CDC OpenData Bereich heruntergeladen werden. Hier werden die Wetterdaten über FTP und HTTPS zum Download angeboten. Unter der URL https://www.dwd.de/DE/leistungen/cdc/cdc_ueberblick-klimadaten. html wird eine gute Übersicht über die zum Download angeboten Daten geboten. Die Werte für die aktuelle Lufttemperatur können über https://opendata.dwd.de/ climate_environment/CDC/observations_germany/climate/10_minutes/air_temperature/ now/ abgerufen werden. Historische Daten können über https://opendata.dwd.de/ climate_environment/CDC/observations_germany/climate/10_minutes/air_temperature/ now/ abgerufen werden.

Aktuell werden auf der Webseite, für die aktuelle Lufttemperatur, ca 480 Dateien zum Download angeboten. Die meisten dieser Dateien entsprechen jeweils einer Messstation und je nach Tageszeit kann deswegen die Menge der Zeilen in der Datei variieren, weil immer um 00:00 eine neue Datei angefangen wird. In den Zip-Dateien finden sich außerdem Metadaten über die Messtationen. Die eigentlichen Daten sind als Commaseparated values (CSV) formatiert und sehen aus wie in Listing 2 gekürzt dargestellt In der CSV Datei gibt es 9 Felder. Der Inhalt der Felder wird in Tabelle 1 beschreiben.

STATIONS_ID;MESS_	DATUN	1; QN;P	P_10;T	Γ_10;TΝ	15_10;R	F_10;TD_10;e	or
73;202205120000;	2;	-999;	ī2.9;	11.2;	84.2;	10.3;eor	
73;202205120010;	2;	-999;	12.7;	11.2;	84.9;	10.2;eor	
73;202205120020;	2;	-999;	12.9;	11.4;	83.0;	10.1;eor	
73;202205120030;	2;	-999;	12.4;	10.7;	86.9;	10.3;eor	
73;202205120040;	2;	-999;	12.4;	10.5;	86.2;	10.2;eor	
73;202205120050;	2;	-999;	12.3;	10.3;	85.5;	9.9;eor	
73;202205120100;	2;	-999;	12.1;	10.1;	88.1;	10.2;eor	
73;202205120110;	2;	-999;	11.7;	9.9;	90.1;	10.1;eor	
73;202205120120;	2;	-999;	11.7;	10.0;	89.0;	10.0;eor	

Listing 2: Wetterdaten CSV Quelle: DWD Wetterdaten CSV Download

In der Datei zehn_now_tu_Beschreibung_Stationen.txt werden die Wetterstationen beschrieben. Diese Datei ist nicht als CSV Datei formatiert sondern als Tabelle und erhält Daten über die Wetterstationen. Die Daten der Stationen in der heruntergeladenen Textdatei stimmen mit den Daten der Hauptamtliches Messnetz Karte überein. Allerdings enthält die Textdatei nicht alle Stationen sondern nur Station für die auch Messwerte im Datensatz hinterlegt sind. Die Bedeutung der einzelnen Spalten der Tabelle sind in der Tabelle 2 beschreiben.

3.7.2 Wetterdaten abrufen

Um die Daten auswerten zu können müssen diese als erstes heruntergeladen und entpackt werden. Dazu wird mithilfe von BeautifulSoup aus der Hypertext Markup Language (HTML) Seite des DWD für jede Datei eine URL ausgelesen. Die so gewonnen URLs können dann mithilfe einer Schleife heruntergeladen werden. Um den Messwerten eine Station zuzuordnen zu können wird als erstes die Datei mit den Station verarbeitet. Für jede Station wird Objekt erstellt und in ein dictionary gespeichert. Dadurch kann in im dictionary einfach über die STATIONS_ID die passende Station gefunden werden. Weil diese Datei allerdings nicht CSV formatiert ist musst die Datei auf eine andere Art ausgewertete werden. Um die einzelnen Felder aus einer Zeile zu bekommen wird immer so lange gelesen bis wieder ein bestimmte Anzahl von Leerzeichen hintereinander erkannt worden ist. Die Zeichen zwischen den Leerzeichen sind dann ein ausgelesenes Feld. Nachdem die Stationsdaten ausgewertet worden sind, werden die CSV Datein in einer Schleife entpackt und mithilfe der Bibliothek Pandas in ein Dataframe umgewandelt. Das so erzeugte Dataframe wir im letzten Schritt mit den Daten der Stations zusammengeführt und als Datenpunkt in InfluxDB geschrieben. Weitere Erklärungen und der Code selbst kann im angehängten Jupyter notebook eingesehen und ausgeführt

Feld Name	Bedeutung
STATION_ID	Gibt an von welcher Station die Werte stammen
MESS_DATUM	Gibt an wann gemessen wurde im Format %Y%m%d%H%M. Also
	Jahr Monat Tag Stunde Minute als eine zusammengeschriebene
	Zahl.
QN	Gibt die Qualität der Messwerte an. Hier gibt es die Werte 1 bis 3
	1. nur formale Kontrolle bei Dekodierung und Import
	2. Kontrolle mit individuell festgelegten Kriterien
	3. ROUTINE automatische Kontrolle und Korrektur mit Soft-
	ware (QUALIMET)
PP_10	Luftdruck auf Stationshöhe
10	Lufttemperatur auf 2m Höhe
TM5_10	Lufttemperatur auf 5cm Höhe
TD_10	relative Luftfeuchtigkeit auf 2m Höhe
eor	END OF RECORD Bedeutet die Zeile ist zu Ende.

 Tabelle 1: Bedeutung der CSV Felder

Quelle: Eigene Darstellung https://opendata.dwd.de/climate_environment/CDC/observations_ germany/climate/10_minutes/air_temperature/now/BESCHREIBUNG_obsgermany_climate_ 10min_tu_now_de.pdf

werden.

3.7.3 Wetterdaten Verarbeiten

Tabelle 2: Bedeutung Stationstabellen Felder

Feld Name	Bedeutung
STATION_ID	Gibt an von welcher Station die Werte stammen
von_datum	Datum seit dem die Station aktiv ist.
bis_datum	Hohe der Station.
Stationshoehe	Hohe über dem Normalnullpunkt
geoBreite	Breitengrad der Station.
geoLaenge _10	Längengrad der Stations.
Stationsname	Name der Station.
Bundesland	Bundesland in dem die Station steht.

Quelle: Vergleich der Werte mit der Hauptamtliches Messnetz Karte https://opendata.dwd.de/ climate_environment/CDC/observations_germany/climate/10_minutes/air_temperature/now/ zehn_now_tu_Beschreibung_Stationen.txt

4 Zusammenfassung

Anhang

Anhangsverzeichnis

Anhang 1:	InfluxDB Python Code	16
Anhang 2:	Times Seires DB Rangliste	17
Anhang 3:	InfluxDB Webinterface Screenshot	18

Anhang 1 InfluxDB Python Code

```
Listing 3: InfluxDB Python Daten schreiben Quelle: Influxdata (2022d)
```

```
import influxdb_client
1
      from influxdb_client.client.write_api import SYNCHRONOUS
2
3
      bucket = "<my-bucket>"
4
      org = "<my-org>"
\mathbf{5}
      token = "<my-token>"
6
      url="http://influxdb:8086"
7
8
      client = influxdb_client.InfluxDBClient(
9
          url=url,
10
          token=token,
11
12
          org=org
      )
13
14
      write_api = client.write_api(write_options=SYNCHRONOUS)
15
16
      p = influxdb_client.Point("my_measurement")
17
          .tag("location", "Paderborn")
18
          .field("temperature", 25.3)
19
      write_api.write(bucket=bucket, org=org, record=p)
20
```

Listing 4: InfluxDB Python Daten lesen Quelle: Influxdata (2022d)

```
query_api = client.query_api()
1
      query = 'from(bucket:"my-bucket")\
2
      |> range(start: -10m)\
3
      > filter(fn:(r) => r._measurement == "my_measurement")\
\mathbf{4}
      |> filter(fn: (r) => r.location == "Prague")\
\mathbf{5}
      > filter(fn:(r) => r._field == "temperature" )'
6
      result = query_api.query(org=org, query=query)
7
      results = []
8
      for table in result:
9
          for record in table.records:
10
             results.append((record.get_field(), record.get_value()))
11
12
      print(results)
13
      [(temperature, 25.3)]
14
```

Anhang 2 Times Seires DB Rangliste

	-		394 Systeme im Ranking, Mai 2022				
	Rang						
Mai 2022	Apr 2022	Mai 2021	DBMS	Datenbankmodell	Mai 2022	Apr 2022	Mai 2021
1.	1.	1.	Oracle 🕂	Relational, Multi-Model 👔	1262,82	+8,00	-7,12
2.	2.	2.	MySQL 🖶	Relational, Multi-Model 👔	1202,10	-2,06	-34,28
3.	3.	3.	Microsoft SQL Server 🔠	Relational, Multi-Model 👔	941,20	+2,74	-51,46
4.	4.	4.	PostgreSQL 🚹 💭	Relational, Multi-Model 👔	615,29	+0,83	+56,04
5.	5.	5.	MongoDB 🚹	Document, Multi-Model 👔	478,24	-5,14	-2,78
6.	6.	个 7.	Redis 🛨	Key-value, Multi-Model 👔	179,02	+1,41	+16,85
7.	1 8.	4 6.	IBM Db2	Relational, Multi-Model 👔	160,32	-0,13	-6,34
8.	4 7.	8.	Elasticsearch 🔁	Suchmaschine, Multi-Model 📷	157,69	-3,14	+2,34
9.	9.	↑ 10.	Microsoft Access	Relational	143,44	+0,66	+28,04
10.	10.	4 9.	SQLite 🚦	Relational	134,73	+1,94	+8,04
11.	11.	11.	Cassandra 🛨	Wide column	118,01	-3,98	+7,08
12.	12.	12.	MariaDB 🛨	Relational, Multi-Model 👔	111,13	+0,81	+14,44
13.	13.	13.	Splunk	Suchmaschine	96,35	+1,11	+4,24
14.	14.	1 27.	Snowflake 🞛	Relational	93,51	+4,06	+63,46
15.	15.	15.	Microsoft Azure SQL Database	Relational, Multi-Model 👔	85,33	-0,45	+14,88
16.	16.	16.	Amazon DynamoDB 🚹	Multi-Model 👔	84,46	+1,55	+14,39
17.	17.	4 14.	Hive 🕂	Relational	81,61	+0,18	+5,42
18.	18.	4 17.	Teradata 🖶	Relational, Multi-Model 👔	68,39	+0,82	-1,59
19.	19.	19.	Neo4j 🔁	Graph	60,14	+0,62	+7,91
20.	20.	20.	Solr	Suchmaschine, Multi-Model 🛐	57,26	-0,48	+6,07
21.	21.	4 18.	SAP HANA 🖶	Relational, Multi-Model 👔	55,09	-0,71	+2,33
22.	22.	22.	FileMaker	Relational	52,27	-0,64	+5,55
23.	1 24.	1 24.	Google BigQuery 🚦	Relational	48,61	+0,63	+10,98
24.			Databricks	Multi-Model 👔	47,85		
25.	4 23.	4 21.	SAP Adaptive Server	Relational, Multi-Model 👔	47,78	-0,58	-2,19
26.	4 25.	4 23.	HBase 🔁	Wide column	43,19	-1,14	-0,05
27.	4 26.	4 25.	Microsoft Azure Cosmos DB 🖶	Multi-Model 👔	40,22	-0,12	+5,51
28.	4 27.	28.	PostGIS	Spatial DBMS, Multi-Model 👔	31,82	-0,23	+1,98
29.	4 28.	29.	InfluxDB 🔂	Time Series, Multi-Model 👔	29,55	-0,47	+2,38
30.	4 29.	4 26.	Couchbase 🚹	Document, Multi-Model 👔	28,38	-0,67	-1,85

Abbildung 1: DB-Engines Ranking

 $\textbf{Quelle: } https://db-engines.com/de/ranking?msclkid{=}4f2a29e5d08811ec95ccd74f8f5146ab$

Anhang 3 InfluxDB Webinterface Screenshot



Quelle: Eigener Screenshot InfluxDB Webinterface

Abbildung 3: InfluxDB Load Data Source



Quelle: Eigener Screenshot InfluxDB Webinterface

Abbildung 4: InfluxDB Load Data Bucket

$\leftarrow \ \rightarrow $	C	O D localhost:8086/o	rgs/427622a51a40f316/loa	ad-data/buck			숪		<u>*</u> *	θ	≡
	Load D	ata									
) Data	Sources	Buckets Teleg	raf Scrapers	API Tok	ens						
Explore	Q Filter buck		Sort by Name (A \rightarrow Z)					H	Create E	Bucket	
Books Boords	_monitor System Bucke	ing et │ Retention: 7 days						Wha Buc	t is a ket?		
Tasks	_tasks System Bucke	et 🕴 Retention: 3 days						A bu nam wher data	cket is a ed locatior e time ser is stored.	n ies All	
Alerts	dwd_now Retention: 1 d + Add a lat	lays ID: c1559afb5e44a2; bel				+ Add Data	Settings	buck Rete dura that point	ets have a ntion Poli tion of tim each data persists.	e	
	Test Retention: Fo + Add a lat	rever ID: e522f42b6cf70				+ Add Data	Settings	Here write your	s how to data into bucket.		
	test-buck Retention: Fo + Add a lat	ced rever ID: 645ec16c035ba bel				+ Add Data	Settings				

Quelle: Eigener Screenshot InfluxDB Webinterface

Create Bucket X						
Name*						
Give your bucket a name						
Delete Data						
Never Older Than						
Cancel Create						

Abbildung 5: InfluxDB Load Data Bucket hinzufügen

Quelle: Eigener Screenshot InfluxDB Webinterface

Abbildung 6: InfluxDB Load Data API Tokens



Quelle: Eigener Screenshot InfluxDB Webinterface

Abbildung 7: InfluxDB Load Data API Token hinzufügen

Generate Read/Write API Token ×									
Description									
	Read		Write						
All Buckets	S	coped	All Bucket	ts S	coped				
BUCKETS	Select All	Deselect All	BUCKETS	Select All	Deselect All				
Test			Test						
dwd_now			dwd_now						
test-bucked			test-bucked						
		X Cancel	✓ Save						

Quelle: Eigener Screenshot InfluxDB Webinterface

Quellenverzeichnis

Monographien

- Silaparasetty, Nikita (2020). Machine Learning Concepts with Python and the Jupyter Notebook Environment. Apress Berkeley CA, S. 91–118.
- Stender, Daniel (2020). Cloud-Infrastrukturen: Infrastructure as a Service So geht moderne IT-Infrastruktur. Das Handbuch für DevOps-Teams und Administratoren. Rheinwerk Computing, S. 54–68.

Internetquellen

- ComputerWeekly.de, Redaktion (2021). Definition Zeitreihendatenbank (Time Series Database, TSDB). URL: https://datascience.eu/wiki/what-the-heck-is-time-series-data-and-why-do-i-need-a-time-series-database/ (besucht am 21. Mai 2021).
- Data-Science-Team (2020). What the heck is time-series data (and why do I need a timeseries database)? URL: https://www.computerweekly.com/de/definition/ Zeitreihendatenbank-Time-Series-Database-TSDB (besucht am 9. Mai 2022).
- Dix, Paul (2021). Why Time Series Matters for Metrics, Real-Time Analytics and Sensor Data. URL: http://get.influxdata.com/rs/972-GDU-533/images/ why%20time%20series.pdf (besucht am 10. Mai 2022).
- Fangman, Sam (2019). The Time Has Come for a New Type of Database. URL: https: //medium.datadriveninvestor.com/the-time-has-come-for-a-new-typeof-database-47cf8df1667a (besucht am 10. Mai 2022).
- Hajek Vlasta Pour Ales, Kudibal Ivan (2019). Why Time Series Matters for Metrics, Real-Time Analytics and Sensor Data. URL: http://get.influxdata.com/rs/ 972-GDU-533/images/why%20time%20series.pdf (besucht am 27. Mai 2022).

- hazelcast (2022). *Time Series Database*. URL: https://hazelcast.com/glossary/ time-series-database/ (besucht am 10. Mai 2022).
- Influxdata (2021). Why You Should Migrate from SQL to NoSQL for Time Series Data. URL: https://www.influxdata.com/from-sql-to-nosql/ (besucht am 20. Mai 2022).
- Influxdata (2022a). API Quick Start. URL: https://docs.influxdata.com/ influxdb/v2.2/api-guide/api_intro/ (besucht am 21. Mai 2022).
- Influxdata (2022b). Line protocol. URL: https://docs.influxdata.com/influxdb/ v2.2/reference/syntax/line-protocol/ (besucht am 29. Mai 2022).
- Influxdata (2022c). Manage API tokens. URL: https://docs.influxdata.com/ influxdb/cloud/security/tokens/#all-access-token (besucht am 29. Mai 2022).
- Influxdata (2022d). Python client library. URL: https://docs.influxdata.com/ influxdb/v2.2/api-guide/client-libraries/python/ (besucht am 18. Juni 2022).
- Influxdata (2022e). Use InfluxDB client libraries. URL: https://docs.influxdata. com/influxdb/v2.2/api-guide/client-libraries/ (besucht am 18. Juni 2022).
- Influxdata (2022f). Write data with the InfluxDB API. URL: https://docs. influxdata.com/influxdb/v2.2/write-data/developer-tools/api/ (besucht am 29. Mai 2022).
- Seeger, Marc (2009). Key-Value stores: a practical overview. URL: http://blog.marcseeger.de/assets/papers/Ultra_Large_Sites_SS09-Seeger_Key_Value_ Stores.pdf (besucht am 19. Juni 2022).
- solid-IT-gmbh (2022). DB-Engines Ranking. URL: https://db-engines.com/de/ ranking (besucht am 10. Mai 2022).

Ehrenwörtliche Erklärung

Hiermit erkläre ich, dass ich die vorliegende Studienarbeit - Entwurf selbständig angefertigt habe. Es wurden nur die in der Arbeit ausdrücklich benannten Quellen und Hilfsmittel benutzt. Wörtlich oder sinngemäß übernommenes Gedankengut habe ich als solches kenntlich gemacht. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Paderborn, 19. Juni 2022

Henrik Mertens