# Praxisnahe Erfahrungen aus dem Datenqualitäts-Dschungel

Florian Gräbe, Marcel Spitzer

#### **Team inovex**

Karlsruhe · Köln · München · Hamburg Berlin · Stuttgart · Pforzheim · Erlangen



# **Marcel Spitzer**





Marcel Spitzer

#### Florian Gräbe





Florian Gräbe

Data & ML Engineer (2016)

"Struggles his way through the data wilderness"

Data & ML Engineer (2022)

"Driven by good Data, ML and Dancing"





Wir freuen uns auf Ihren Besuch an **Stand 2** im Foyer!

# WE ARE INOVEX – INNOVATE, INTEGRATE, EXCEED.

Ihr Partner für die digitale Transformation mit dem Leistungsspektrum

Data & AI
Application Development
Skalierbare IT-Infrastrukturen
Training & Coaching



Ausschließlich festangestellte Mitarbeitende



Wachstum durch Innovationen und erfolgreiche Projekte





**SOFTWARE · DATA & AI · INFRASTRUCTURE** 

# Begeben wir uns in die Wildnis

- 1. Sind wir im Daten-Dschungel?
- 2. Risiken und Gefahren
- 3. Sicheres Fortbewegen
- 4. Survival Ausrüstung





**Undurchsichtig** 

#### Verwachsen



Unterschiedliche Blickwinkel

**Chaotisch / Wild** 



#### Verwachsen

#### **Undurchsichtig**

Welche Regeln und Gesetze gelten (GDPR, ...)?

Woher kommen die Daten?



Unterschiedliche Blickwinkel

**Chaotisch / Wild** 



Wozu sind die Daten nutzbar?

Kann ich den Daten vertrauen?

#### **Undurchsichtig**

Welche Regeln und Gesetze gelten (GDPR, ...)?

Woher kommen die Daten?

#### Verwachsen



Unterschiedliche Blickwinkel

**Chaotisch / Wild** 



Wozu sind die Daten nutzbar?

Kann ich den Daten vertrauen?

#### **Undurchsichtig**

Welche Regeln und Gesetze gelten (GDPR, ...)?

Woher kommen die Daten?

#### Verwachsen



**Chaotisch / Wild** 

Wie kann ich auf die Daten zugreifen?

Wer besitzt welche Daten?

Unterschiedliche Blickwinkel



Wozu sind die Daten nutzbar?

Kann ich den Daten vertrauen?

#### **Undurchsichtig**

Welche Regeln und Gesetze gelten (GDPR, ...)?

Woher kommen die Daten?

#### Verwachsen



**Chaotisch / Wild** 

Wie kann ich auf die Daten zugreifen?

Wer besitzt welche Daten?

Unterschiedliche Blickwinkel

In welcher Qualität und Formaten liegen die Daten vor?



### Was verstehen wir unter Datenqualität?

- Allgemein: der Grad der Übereinstimmung von Daten mit der Realität
- in der Praxis jedoch oft schwierig bis unmöglich auf diese Weise zu messen
- Deshalb: Abweichung von zuvor festgelegten Annahmen
- Annahmen ergeben sich aus der jeweiligen Domäne, deshalb regelmäßiger Austausch mit Experten erforderlich



# **Gefahren im Daten-Dschungel**

- Auf Basis von fehlerhaften Daten können falsche Schlussfolgerungen gezogen und suboptimale Entscheidungen getroffen werden
- Die Qualität eines AI- oder ML-Modells hängt maßgeblich von den Trainingsdaten ab (Garbage-in-Garbage-out)
- Beispielszenarien, in denen fehlerhafte Daten großen Schaden anrichten können:
  - o Empfehlungssysteme
  - Absatzprognose
  - Investitionsentscheidungen





# Mit Plan sicher durch den Daten-Dschungel bewegen

- Kooperatives Erarbeiten von Annahmen
- 2. Definition von **Strategien** für den Umgang mit Auffälligkeiten und Fehlern
- 3. **Integration** von Datenvalidierung in Datenpipelines
- 4. **Überwachen** der Validierungsergebnisse und aktualisieren der Annahmen



### Herausforderung: Anforderung sammeln & konsolidieren





# Datenqualität umfasst viele Dimensionen

Datentypen Regex Pattern Vollständigkeit **Nullability** Schwellenwerte syntaktische semantische Korrektheit Korrektheit Dates/ Wertebereiche Timestamps Korrelationen Stakeholder **Aktualität** Relevanz Bedeutung

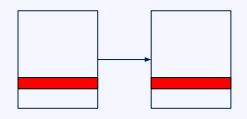


# Kritikalität bestimmt den Umgang mit fehlerhaften Daten

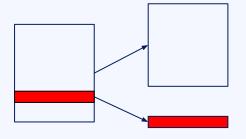
- Je nach **Datensatz** und **Annahme** sind unterschiedliche Reaktionen erforderlich
- Unabhängig von der Reaktion sind
   Benachrichtigungen in jedem Falle sinnvoll
- Wurden fehlerhafte Daten identifiziert gibt es verschiedene Möglichkeiten, darauf zu reagieren:



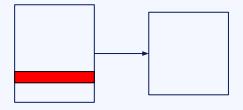
# Kritikalität bestimmt den Umgang mit fehlerhaften Daten



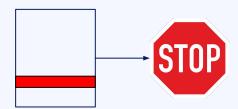
(a) weiterverarbeiten



(c) aussortieren

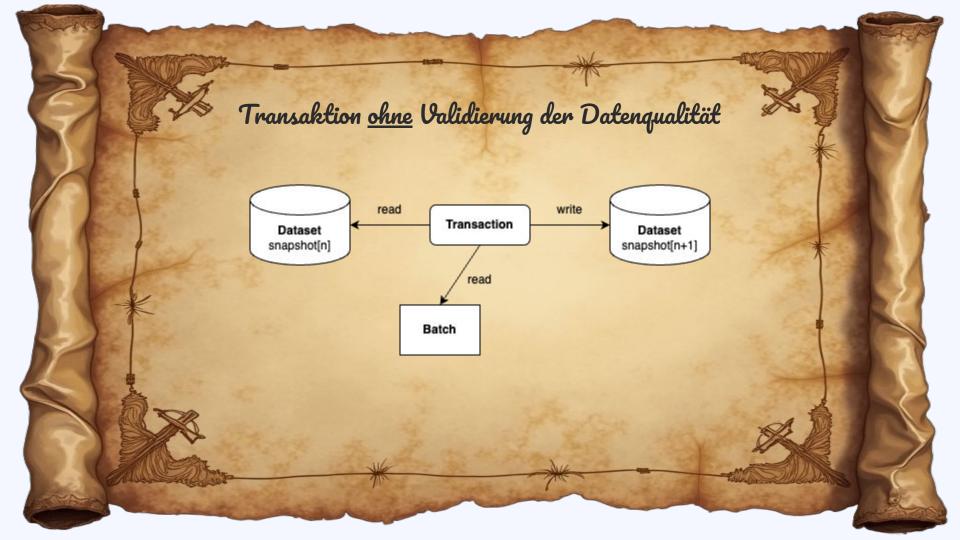


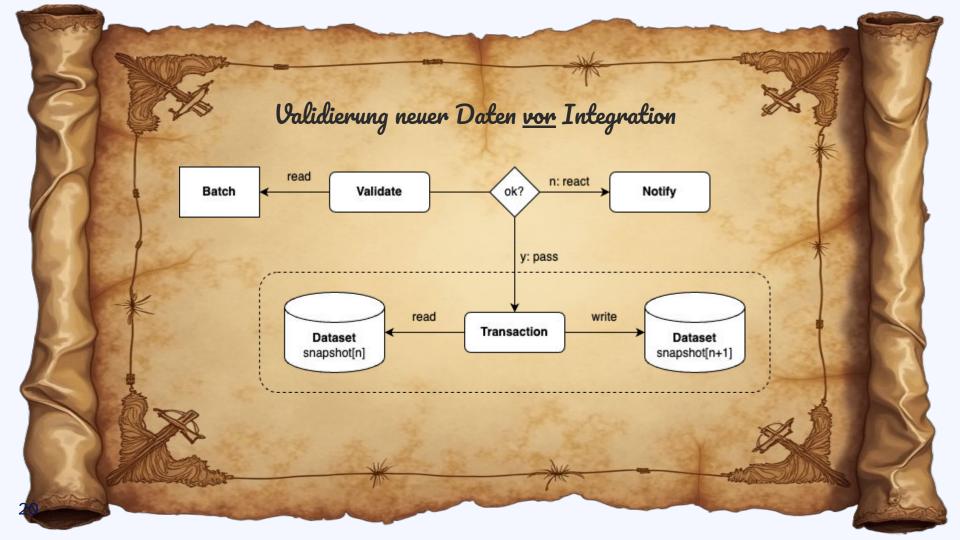
(b) ignorieren

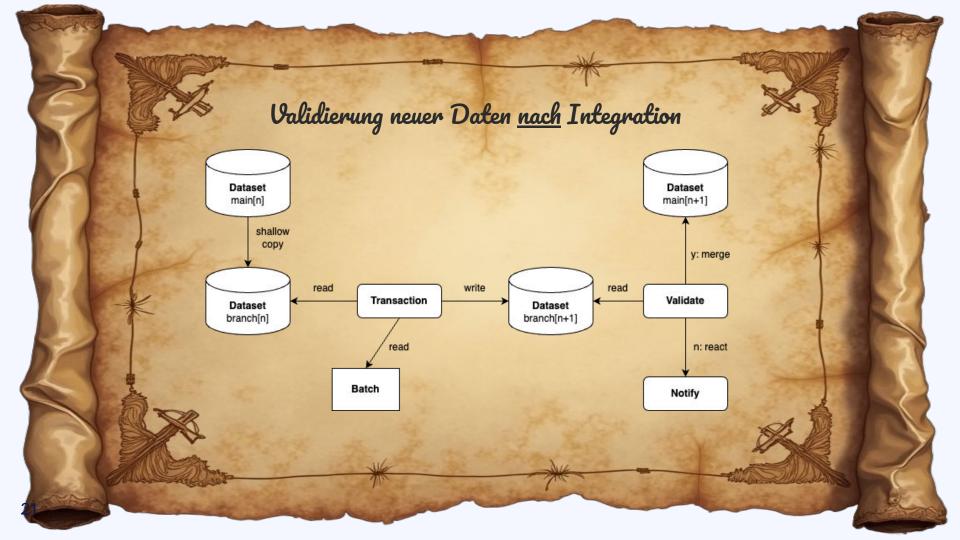


(d) stoppen



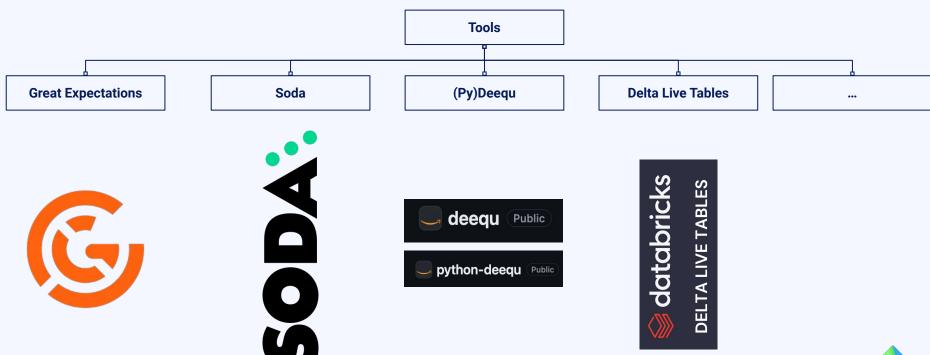








# Ein breites Angebot an Werkzeugen





# Great Expectations: Die etablierte Toolbox für den Ernstfall



"Great Expectations is the leading tool for validating, documenting, and profiling your data to maintain quality and improve communication between teams"

- "Eierlegende Wollmilchsau" im Datenqualitäts-Dschungel
- Breites Einsatzspektrum, aber steile Lernkurve
- Umfangreiches Set an **Expectation Optionen**
- Große Community und Verbreitung



#### **Great Expectations: Context Setup**

```
Data Context
              Data Source
                           expectation.json
                                             Checkpoint
                                                             Data Docs
from great_expectations.data_context import EphemeralDataContext,
FileDataContext
# Fluechtiger in Memory Context
context = EphemeralDataContext(project_config=project_config)
# Context mit Filesystem Backend
path_to_folder = "./"
context = FileDataContext.create(project_root_dir=path_to_folder)
```



#### **Great Expectations: Datenquellen**

```
Data Context
               Data Source
                             expectation.json
                                                Checkpoint
                                                                 Data Docs
# Nutzen einer Postgres als Datenquelle
pg datasource = context.sources.add postgres(
    name="pg datasource", connection string=PG CONNECTION STRING
# Nutzen einer Datei in S3 als Datenquelle
data source = context.sources.add pandas s3(
   name=source name, bucket=bucket name
data asset = data source.add csv asset(
   name=asset_name, batching_regex=batching_regex, s3_prefix=s3_prefix
data asset.build batch request()
```

< 0.18.x

# **Great Expectations: Expectation Suite**

```
Data Context
               Data Source
                             expectation.json
                                                 Checkpoint
                                                                   Data Docs
 "expectation_suite_name": "d2d_expectations",
 "ge_cloud_id": null,
 "expectations": [
 "expectation_type": "expect_column_values_to_be_between",
 "kwargs": {
   "column": "listener_count",
   "max_value": 1000000,
   "min_value": 1,
 "meta": {}
},
```



<0.18.x

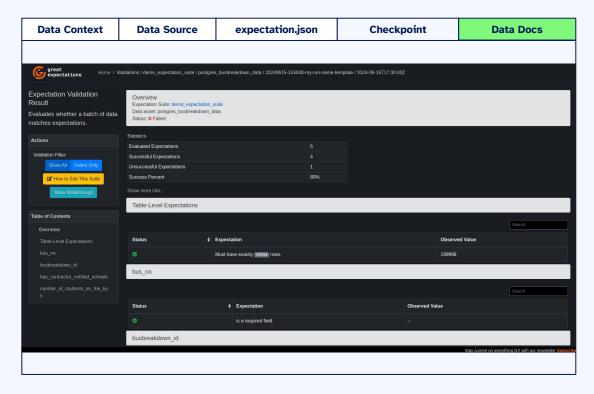
#### **Great Expectations: Checkpoint**

```
Data Context
               Data Source
                             expectation.json
                                                 Checkpoint
                                                                  Data Docs
validations = [
       "batch_request": batch_request,
       "expectation_suite_name": expectation_suite,
       "action_list": action_list,
checkpoint = context.add_or_update_checkpoint(
   name="Checkpoint",
   validations=validations,
   runtime_configuration={}},
checkpoint.run(run name=run name)
```



<0.18.x

### **Data Docs: Ergebnisse als HTML**





#### Der Challenger im DQ Dickicht

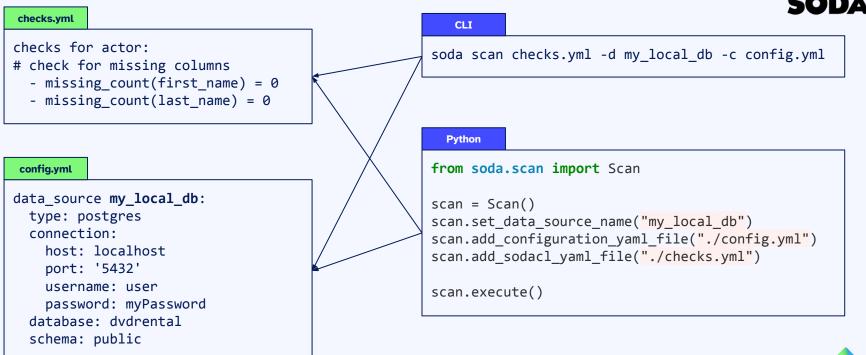


"Build reliable data products and pipelines with Soda, the GenAI-first platform for data quality. Embed tests into your workflows and monitor data quality health any way you like-through out-of-the-box observability or declarative testing."

- Breites Angebot an Checks und Datenquellen
- Möglichkeit zur Erstellung individueller Checks
- **Intuitives** Nutzungskonzept
- aktuell weniger verbreitet



#### Datasets und Checks werden in Soda in YAML definiert





#### Datasets und Checks werden in Soda in YAML definiert



```
CLI
soda scan checks.yml -d my_local_db -c config.yml
              Scan summary:
              2/2 checks PASSED:
                  actor in my_local_db
                    missing_count(last_name) = 0 [PASSED]
                    missing_count(first_name) = 0 [PASSED]
              All is good. No failures. No warnings. No errors.
```



#### Soda ist ein starker Konkurrent zum etablierten GX

	SODA:	Great Expectations
Integrationsmöglichkeiten	+	+
Funktionsumfang	+	+
Nutzungskonzept	+	-
Dokumentation	+	-
Verbreitung/Community	-	+



# Fazit: Überleben im Daten-Dschungel

#### Risiken und Gefahren kennen:

- Konsequenzen fehlerhafter Daten
- Vorsorge ist besser als Nachsorge

#### Verbündete mit einbeziehen:

Domänenexperten, Nutzer, Stakeholder, ...

#### **Stets den Kompass im Blick behalten:**

 Iterative Herangehensweise, inkrementelle Verbesserung

#### Geschulter Umgang mit Werkzeugen:

- Integrierte Datenvalidierung
- Strategien f
  ür den Umgang mit Auff
  älligkeiten



# Vielen Dank!



**Marcel Spitzer Data & ML Engineer** 

marcel.spitzer@inovex.de



Florian Gräbe **Data & ML Engineer** 

florian.graebe@inovex.de





# Vortragsfolien zum Download

Besuchen Sie uns gerne auch an unserem Stand im Foyer!



