

Data Quality and Validation in ML Pipelines

Great Expectations, Deequ, and TensorFlow Data Validation

By: Sunil Kumar Mudusu

Why Data Quality Matters



Models Need Good Data

ML models need good data to learn and work well



Bad Data = Bad Predictions

Errors, wasted time, broken models



Can't Fix With Better Models

Even small issues in data can break your model





Common Data Problems

Missing Values

Empty fields or null values

Wrong Data Types

Numbers stored as text

Schema Drift

Structure changes over time

Outliers

Strange values breaking patterns



What is Data Validation?



Verify Data Quality

Check data is clean, complete, correct



Early Detection

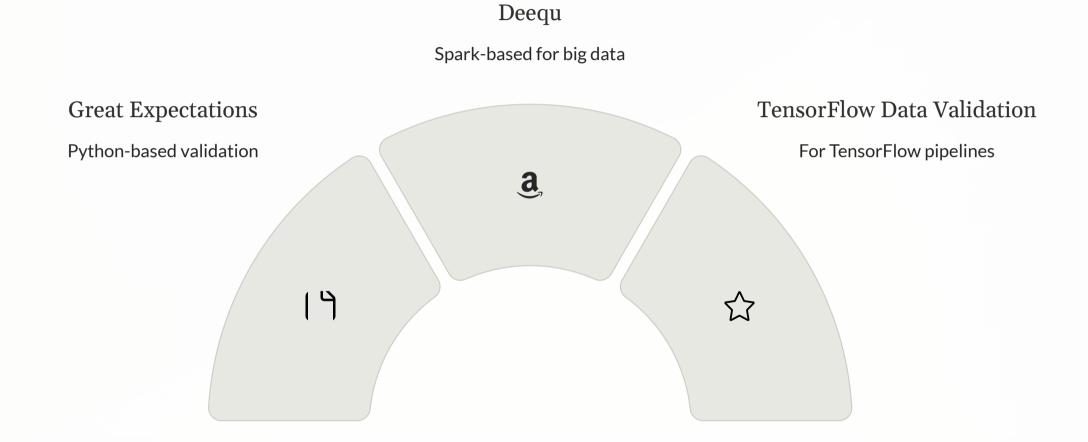
Catch problems before model training

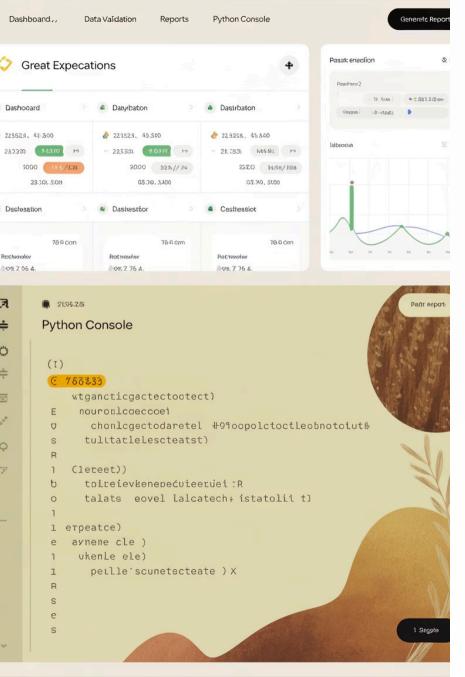


Automation

Can be automated in every ML pipeline

Meet the Tools





Great Expectations



Python-Based

Familiar for data scientists

Easy Rules

Simple data expectations

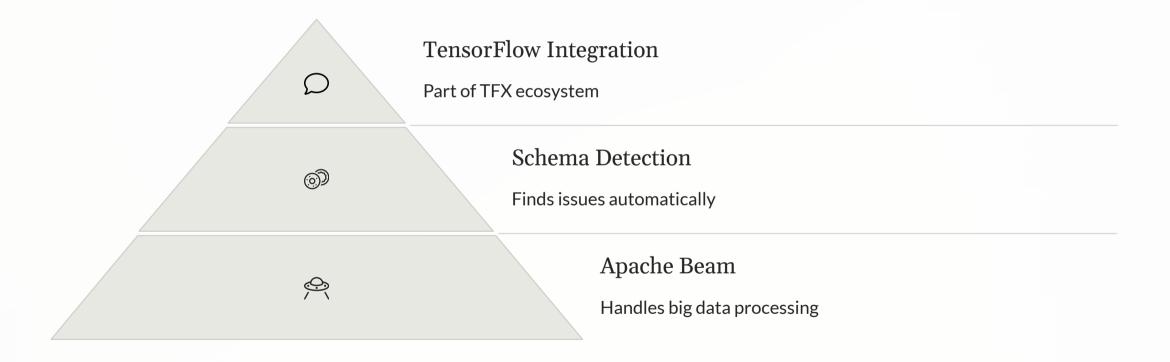
Clear Reports

Human-readable results

Deequ

Amazon-Built pay Runs on Apache Spark Big Data Ready Handles large datasets efficiently **Quality Metrics** 000 Tracks data quality over time

TensorFlow Data Validation



Comparing the Tools

| Tool | Language | Best Use |
|--------------------|--------------|----------------------------|
| Great Expectations | Python | Easy validation, reporting |
| Deequ | Scala/Python | Big data, streaming |
| TFDV | Python | ML pipelines |



Choosing the Right Tool

Assess Your Needs

Consider data size, tools, team skills

Match Tool to Workflow

Great Expectations: readable reports

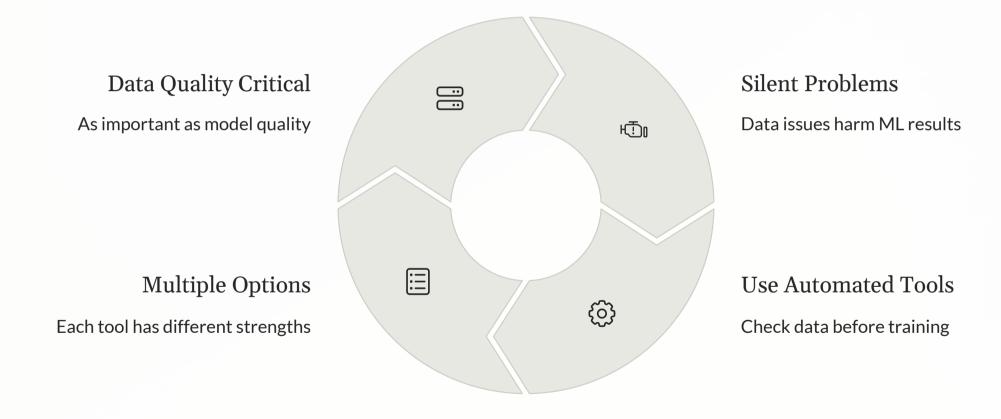
Deequ: large-scale data

TFDV: TensorFlow pipelines

Implement and Iterate

Start simple, expand validation over time

Key Takeaways





Conclusion



Good Data

Leads to better models



Right Tool

Choose for your workflow

Thank You