

# Einführung in MLOps

## 11 PROCESSING- UND PREDICTION MODES, FEATURES

Tobias Mérinat

teaching2025@fsck.ch

Lucerne University of  
Applied Sciences and Arts

**HOCHSCHULE  
LUZERN**

DEPARTMENT OF INFORMATION TECHNOLOGY  
Lucerne University of Applied Sciences and Arts  
6343 Rotkreuz, Switzerland

14. und 15. Februar 2025

Drei grundlegende Arten, wie Daten verarbeitet werden

- Request-Response
- Batch Processing
- Stream Processing

- stateless
- synchron
- einfach skalierbar

- stateful
- bounded
- parallelisierbar

- stateful
- unbounded
- komplex

- Batching und Streaming gemacht für stateful, high Volume Processing
- Request-Response
  - eher für User-Facing Apps verwendet
  - für Data Processing meist nicht geeignet, da im Normalfall stateless
  - kann aber auch die Resultate aus Batch-Processing nachladen und verwenden

- **Features:** Die Inputs eines ML Modells
- **Feature-Engineering:** Berechnung der Features aus Rohdaten
- **Unstrukturierte Daten:** (eher weniger Feature Engineering)
- **Strukturierte Daten:** (eher mehr Feature Engineering)

- Batch-Features
- Streaming-Features
  - Real-Time (RT) Features
  - Near-Real-Time (NRT) Features



## Feature-Arten 2/2

Feature-Typ	Berechnung	Compute Engine			
		(Bsp.)	Latenz	Vorteile	Nachteile
Batch	Vorberechnet in Batch-Prozess	Spark	Stunden bis Tage	einfach, parallelisierbar	keine aktuellen Features
RT	Zum Zeitpunkt der Vorhersage	Python, SQL	<1s	einfach, aktuelle Features	skaliert nicht
NRT	Vorberechnet in Streaming-Prozess	Flink, Quix	Sek. bis Min.	Aktuelle Features, skalierbar	schwierig umsetzbar

- Setzt nur eine (häufig bereits bestehende) Batch-Processing Infrastruktur voraus
- Vorteile
  - Wenig komplex
  - Parallelisierbar
- Nachteile
  - Begrenzte Aktualität
  - u.U unnötige Berechnungen

- Direkt im Rahmen jeder Anfrage berechnet
- Vorteil
  - Top-Aktuell
  - Einfach umsetzbar
- Nachteil
  - Berechnungszeit erhöht direkt Latenz, skaliert deshalb schlecht

- Werden (im Gegensatz zu RT) asynchron berechnet
- Latenz wird nur durch Lookup-Zeit erhöht
- Unterschied zu Batch:
  - Tendentiell häufiger berechnet
  - Mittels Stream-Processing berechnet

Wir unterscheiden verschiedene Stufen

- 1 Offline-Prediction (Batch-Prediction)
- 2 Online-Prediction
  - 1 ausschliesslich mit Batch-Features
  - 2 mit Streaming- und Batch Features
- 3 Continual Learning

# Offline-Prediction (Batch-Prediction)

- Berechnungen erfolgen periodisch und im Voraus (vor der Verwendung)
- Typische Anwendungsfälle:
  - Lead Scoring
  - Demand Forecasting
  - Inventory Management
- Voraussetzungen sind eine Batch-Infrastruktur
- Model Registry wäre eine mögliche Ausbaustufe

- Online Prediction: Wenn die Vorhersage erst gemacht wird, wenn sie benötigt wird
- Nur vorausberechnete Features werden verwendet
- Wenn Batch-Features (zeitlich) *akkurat genug* sind
- Neben der Batch-Infrastruktur wird schneller Key Value Store benötigt

- RT, NRT und Batch Features kombiniert
- Voraussetzung ist eine Stream Processing Infrastruktur



- Fast immer ändern sich Umwelt und somit Input-Daten für Modell über die Zeit
- Dann müssen Modelle periodisch neu trainiert werden
- Dazu werden nicht nur Inputdaten (X) sondern auch das Target (y) benötigt
- Manchmal fällt das Target automatisch im Prozess an (sog. Natural Labels)
- Oft muss es aber manuell erarbeitet werden

- Manuelles stateless Retraining
- Automatisiertes stateless Retraining
- Automatisiertes stateful Training
- Continual Learning

- Retraining, wenn sich der Bedarf zeigt
- Labelling manuell
- Zusammenstellung der Trainingsdaten manuell
- Mit alten und neuen Daten oder nur mit neuen
- Stateless: Modell wird von Grund auf neu trainiert

- Möglich mit Natural Labels
- Weniger fehleranfällig
- Komplex, wenn Inputdaten und Labels zeitlich stark versetzt anfallen

- Auch *inkrementelles Lernen*
- Modelle werden mit aktuellen Daten *weitertrainiert*
- Setzt *Model Lineage* voraus

- Vollkommene Automatisierung
- Nimmt die Performance zu stark ab, wird automatisch neu trainiert
- Preisfragen:
  - Wann sollte neu trainiert werden
  - Ist das neue Modell tatsächlich besser