

Das Data-Leakage-Problem für die prädiktive Instandhaltung

Martin Patrick Pauli, Martin Golz

Hochschule Schmalkalden, Fakultät Informatik, Blechhammer 4, 98573 Schmalkalden

Einleitung

- Reinstwasser-Aufbereitungsanlage der Mikroelektronik-Industrie
- Ziel der Anlage:
 - Ultrareines Wasser mithilfe Umkehrosmose-Filteranlage

Problem

- Unvorhersehbare Zustandsverschlechterung von Filtereinsätzen innerhalb weniger Stunden
- Filterausfall wird nicht immer rechtzeitig erkannt; Produktionsausfall droht

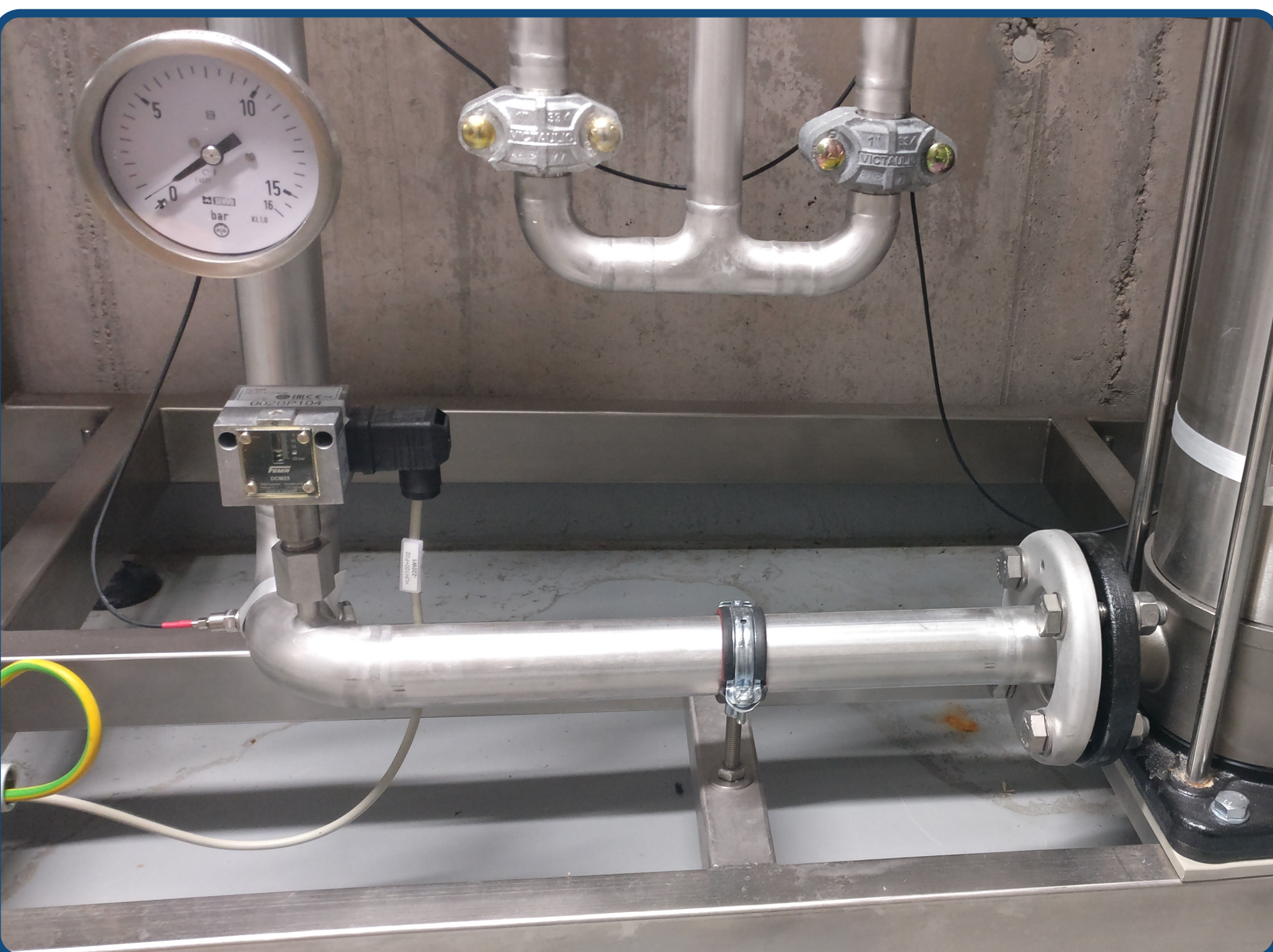
Ziel

- Nicht-invasive Zustandsüberwachung von Filtern
- Frühzeitige Erkennung degenerierter Filter
- Planung des Filterwechsels



Aufbau

- Reinstwasser-Aufbereitungsanlage:
 - Mehrere Filterstufen
- Hauptaggregat:
 - Umkehrosmose-Filter
- Körperschall breitet sich mit geringer Dämpfung in Edelstahl-Rohren aus
- Erfassung der Strömungsgeräusche:
 - Zwei vibroakustische Sensoren
- Sensor 1:
 - Gehäuse des Umkehrosmose-Filter
- Sensor 2:
 - Ausgangsrohr der Vorfilter-Stufe



Daten

- Vibroakustische Signale
 - Körperschall-Zeitreihen
 - Abtastrate: 16.000 Hz
- Zeitraum:
 - Vor Filterwechsel : 48-Stunden-Intervall
 - Nach Filterwechsel: 48-Stunden-Intervall
- Klassen:
 - Degenerierter Filterzustand: vor Filterwechsel
 - Intakter Filterzustand: nach Filterwechsel
- Segmentierung:
 - Nicht überlappend
 - Segmentlänge: 60 s
- Datenumfang:
 - 5 Filterwechsel, 2 Klassen
 - 2.880 Segmente pro Erfassungsintervall
 - Summe: 28.800 Segmente

CI-Methode

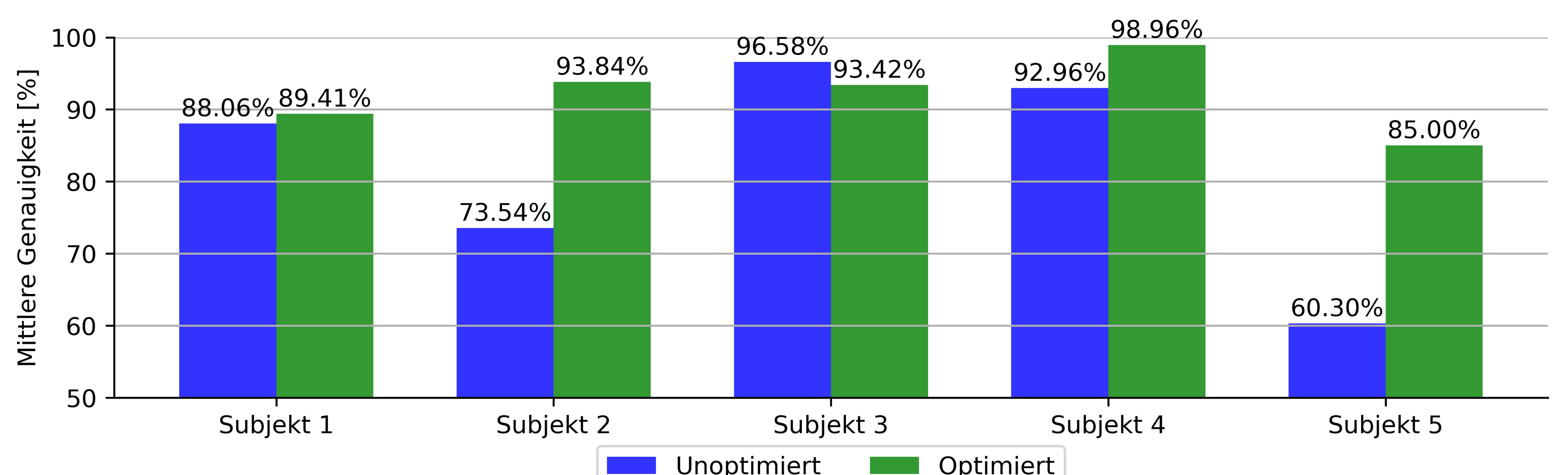
- Gradient Boosting (LightGBM von Microsoft)
- Basis:
 - Ensemble von Entscheidungsbäumen

Kreuzvalidierungsmethoden

- Standard: Repeated random subsampling (RRSS)
 - 80% Trainings-, 20% Validierungsmenge
- Leave-one-subject-out (LOSO)
 - Subjekt: Filterwechsel
 - Je Subjekt:
 - Daten von 2 Filtereinsätzen (vor/nach Wechsel)
 - Strikte Trennung der Subjekte:
 - Einteilung in Trainings- und Validierungsmenge
 - Alle Daten eines Filtereinsatzes:
 - Entweder in Trainings- oder in Validierungsmenge

Ergebnis

- Standard-Kreuzvalidierung (RRSS)
 - 50 Wiederholungen
 - **Klassifikationsgenauigkeit: $99,95 \pm 0,01$ %**
- Kreuzvalidierung ohne Data-Leakage (LOSO)
 - 5 Subjekte = 5 Filterwechsel
 - Balancierung: Gleicher Mengenumfang pro Subjekt
 - Trainingsmenge: 4 Subjekte (80%)
 - Validierungsmenge: 1 Subjekt (20%)
 - **Klassifikationsgenauigkeit: $85,3 \pm 8,7$ %**

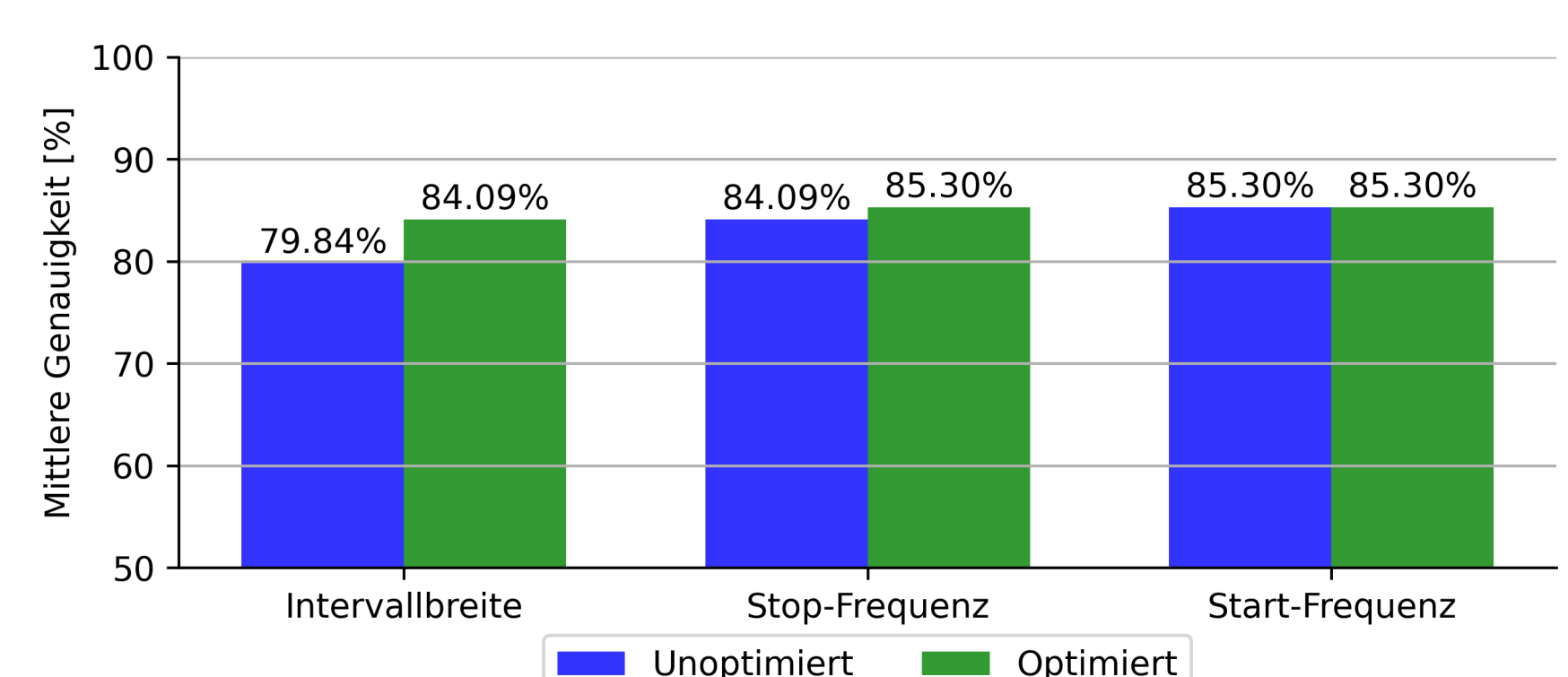


Merkmalsextraktion

- Merkmale:
 - Perzentile der Amplitudenverteilung
 - Momentanfrequenz
 - Entropie-Maße und fraktale Dimensionen
 - Entropie-Maße über DWT-Koeffizientenfolgen
 - Logarithmierte spektrale Leistungsdichten
- Skalierung:
 - Unskaliert
 - Minimum-Maximum-Skalierung
 - Z-Skalierung
 - Quantil-Transformation

Optimierung

- Erste Analysen:
 - LogPSD-Merkmale führen zur maximalen mittleren Genauigkeit
- Weitere Analysen:
 - LogPSD als alleiniges Merkmal
 - Optimierung der Parameter mit RRSS
- Empirische Optimierung:
 - Untere Grenzfrequenz (Start-Frequenz)
 - Obere Grenzfrequenz (Stop-Frequenz)
 - Intervallbreite (Schrittweite)
- Skalierungen: alle wurden untersucht
- Unskalierte LogPSD-Merkmale führen zu maximaler mittlerer Genauigkeit



Diskussion

- Hohe Unterschiede von RRSS im Vergleich mit LOSO weisen auf ein Data-Leakage-Problem hin
- RRSS:
 - Nutzung aller Daten, ohne Beachtung von Subjekten führt zu Data Leakage
 - Klassifikationsgenauigkeiten optimistisch verzerrt
- LOSO:
 - Simuliert Verwendung zukünftiger unbekannter Daten
 - Modell wird unabhängig von einzelnen Subjekten