

# AI-Based Detection of Hedge Fund Fraud

## Section 8 – Reproducibility Statement

Joerg Osterrieder

Zurich University of Applied Sciences (ZHAW)

2025

1. Literature Search Protocol (SALSA Methodology)
2. Search Funnel: 500 → 120 → 105
3. Search Funnel Chart
4. Data Sources: Commercial Databases and Public Filings
5. Database Coverage Chart
6. Limitations: Proprietary Data and Licensing Constraints
7. Code and Implementation References
8. OP1 as Critical Enabler for Reproducibility

## Methodology

Systematic search protocol informed by the **SALSA** (Search, Appraisal, Synthesis, and Analysis) methodology (Grant & Booth, 2009).

### Five Scholarly Databases

1. Scopus
2. Web of Science
3. IEEE Xplore
4. Social Science Research Network (SSRN)
5. Google Scholar

**Search Period:** 2000–2025

Additional: snowballing via reference lists of highly cited papers.

### Structured Boolean Query

- **Domain:** “hedge fund”, “alternative investment”, “investment fund”
- **Fraud:** “fraud”, “manipulation”, “anomaly”
- **Method:** “machine learning”, “artificial intelligence”, “deep learning”, “neural network”
- Searched across titles, abstracts, and keywords

---

Source: Grant & Booth (2009); paper Reproducibility Statement

1. **Initial database searches:** ~500 potentially relevant publications
2. **Title and abstract screening:** reduced to 120 papers warranting full-text review
3. **Full-text review** against inclusion criteria:
  - Peer-reviewed or widely cited preprints
  - Addresses fraud/anomaly detection in investment funds
  - Uses AI/ML methods⇒ **80 papers** met all criteria (core corpus)
4. **Contextual background:** +25 papers on financial regulation, general fraud detection, foundational ML
5. **Total systematically identified:** **105 papers**

## Qualitative Focus

The synthesis focuses on **methodological insights and contextual applicability** rather than exhaustive enumeration.

## Inclusion Criteria

- Peer-reviewed publications or widely cited preprints
- Addresses fraud or anomaly detection
- Focus on investment funds
- Uses AI/ML methods

Source: Paper Reproducibility Statement

ng literature screening stages: 500 initial results → 120 after title/abstract screening → 80 core papers after full-text review → 105 to

- Core corpus of 80 papers supplemented by 25 foundational/contextual references
- Full bibliography includes additional references cited for technical context

---

**Source: Paper Reproducibility Statement**

## Commercial Hedge Fund Databases

- **Lipper TASS** (now Refinitiv)
- **Hedge Fund Research (HFR)**
- **BarclayHedge**
- **Morningstar**

Access requirements:

- Institutional subscriptions
- Licensing restrictions prohibit redistribution
- Widely accessible via university libraries and commercial licenses

## Public Regulatory Filings

- **SEC EDGAR System:**
  - Form ADV (adviser registration)
  - Form D (private offering)
  - Form 13F (quarterly holdings)
- Publicly accessible: [sec.gov/edgar](http://sec.gov/edgar)
- No restrictions on use

## European Regulatory Data

- AIFMD data via national competent authorities
- Access procedures vary by jurisdiction

---

Source: Paper Reproducibility Statement

TASS, HFR, BarclayHedge, Morningstar, SEC EDGAR, AIFMD) mapped against accessibility level (public, institutional subscription, r

- sec *EDGAR is the only fully public source with no access restrictions*
- Commercial databases require institutional subscriptions but are widely available in academia
- Enforcement-labeled datasets are held by regulators and **not publicly accessible**

---

Source: Paper Reproducibility Statement

## Limitations: Proprietary Data and Licensing Constraints

### Proprietary Enforcement Datasets

- Several cited studies use **proprietary enforcement datasets** maintained by regulatory agencies
- Cannot be independently verified due to confidentiality restrictions
- We rely on published methodological descriptions
- Noted where data access limitations prevent full reproducibility

### Alternative Data Licensing

- Some referenced alternative data sources:
  - Satellite imagery
  - Geolocation analytics
  - Sentiment data
- Available only through **vendor-specific licenses**
- Cannot be redistributed
- We prioritize discussion of methods evaluable on **publicly accessible data**

### Implication

Full end-to-end reproducibility of hedge fund fraud detection research remains constrained by **data access limitations endemic to the domain**.

Source: Paper Reproducibility Statement

## Survey Paper: No Original Code

As a qualitative survey, no original computational experiments were conducted and no original code is provided.

## Publicly Available Implementations of Reviewed Methods

Method Category	Library / Framework
Classical ML (RF, GBM, SVM)	scikit-learn (Python), caret (R), XGBoost
Deep learning (FFN, LSTM, autoencoder)	TensorFlow, PyTorch, Keras
Graph neural networks	PyTorch Geometric, Deep Graph Library
NLP / transformers	Hugging Face Transformers, spaCy
Explainability (SHAP, LIME)	shap (Python), lime (Python/R)

- Where specific papers have released code, associated repositories are cited in the paper
- All canonical methods use **mature, well-documented** open-source ecosystems

---

Source: Paper Reproducibility Statement

## The Fundamental Bottleneck

Until public benchmark datasets exist (OP1), full end-to-end reproducibility of hedge fund fraud detection research will remain **constrained by data access limitations endemic to the domain**.

### What OP1 Would Unlock

- **Reproducible comparisons** across studies
- Standardized evaluation protocols
- Cross-method benchmarking on common data
- Progress measurement for the field
- **Accelerated methodological advancement** (as demonstrated by benchmarks in credit card fraud, image classification, NLP)

### Two-Track Path Forward

1. **Synthetic benchmarks**: regime-switching models with injected fraud patterns, calibrated to real distributions
2. **Anonymized regulatory data**: differential privacy protocols enabling regulators to release sanitized enforcement-labeled datasets

Researchers are encouraged to consult the full research agenda (Section 6) for detailed approaches.

---

Source: Paper Reproducibility Statement; Section 6.1.1 (OP1)