# From Models to Decisions: What It Takes to Make Machine Learning Work in Fraud and Credit Risk

## INTRODUCTION:

Machine Learning systems now sit at the core of high-stakes decision-making in fraud detection, credit risk assessment, and customer intelligence. Gradient boosting models, random forests, and deep learning architectures are routinely used to score transactions in milliseconds, evaluate borrower risk at scale, and personalise customer interactions across digital channels. More recently, Generative AI and transformer-based models have expanded these capabilities through embeddings, representation learning, and natural language interfaces, which are layered on top of traditional risk pipelines. Yet despite increasingly sophisticated architectures and higher offline validation metrics, many organisations continue to struggle with rising fraud losses, unstable credit portfolios, and inconsistent model performance in production environments.

The challenge lies not in model selection or algorithmic power, but in how these systems are engineered, evaluated, and governed under real-world constraints. Fraud and credit risk models operate in adversarial, non-stationary settings where concept drift, data leakage, and feedback loops are the norm rather than the exception. Generative AI compounds this complexity by introducing probabilistic outputs, prompt sensitivity, and new sources of model risk. Optimising for AUC or precision-recall in isolation often obscures the real objective: minimising expected loss, controlling tail risk, and ensuring regulatory explainability. To generate durable business value, machine learning in risk-critical domains must be treated as. a continuously monitored decision system—one that tightly integrates data pipelines, model lifecycle management, human oversight, and measurable economic outcomes.

## FEATURE ENGINEERING STILL BEATS ALGORITHMS

There is a reason experienced machine learning practitioners spend far more time on feature engineering than on algorithm selection—especially in fraud detection and credit risk modelling. In widely used credit card fraud datasets, the raw feature space typically consists of anonymised PCA components that are from V1-V8, a transaction timestamp (time), and transaction amount (Amount), alongside a highly imbalanced target label (class). While advanced algorithms such as Gradient Boosting, XGBoost, Random Forests, and Neural Networks can model complex nonlinear relationships, they are fundamentally limited by the representational power of these raw inputs. In practice, the largest performance gains rarely come from switching algorithms; they come from transforming these low-level variables into features that encode behaviour, risk, and temporal context.

### DOMAIN-DRIVEN FEATURE DESIGN USING TRANSACTIONAL VARIABLES

Even when transactional features are anonymised (as with V1-V8), meaningful signals can be extracted by anchoring feature design to domain intuition. For example, raw Amount values are often poorly behaved,

exhibiting heavy tails and extreme skew. A simple log transformation stabilises variance and improves model learning:

```
[10]    # Select Time and Amount
✓ 0s    cols = df[["Time", "Amount"]]

        # Calculate descriptive statistics
        stats = cols.describe().T
        stats
```

| | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| Time | 284807.0 | 94813.859575 | 47488.145955 | 0.0 | 54201.5 | 84692.0 | 139320.500 | 172792.00 |
| Amount | 284807.0 | 88.349619 | 250.120109 | 0.0 | 5.6 | 22.0 | 77.165 | 25691.16 |

Similarly, while the PCA components (V1-V8) are individually opaque, combinations of these features—such as absolute values, interaction terms, or anomaly scores—often carry stronger fraud signals than raw values alone. In practice, gradient-boosted trees trained on engineered transformations of these variables consistently outperform deep neural networks trained directly on the raw components.

## Temporal and Velocity-Based Aggregation

Fraud is inherently temporal. A single transaction rarely signals fraud in isolation; patterns emerge over sequence. The Time variable in the datasets, which represents seconds elapsed since the first transaction, enables powerful velocity and rolling-window features. For example, transaction frequency within a short time horizon is one of the strongest fraud indicators:

```
[12]    #Sort by time
✓ 0s    df= df.sort_values("Time")
        #Rolling transaction count over last 1 hour
        #df["rolling_count"]= df["Amount"].rolling("1H").count()
        df['txn_count_1h']=(
            df.rolling(window=3600, on="Time")['Amount']
            .count()
            .fillna(0)
        )
```

Similarly, rolling monetary aggregates capture sudden spending spikes:

```
[14]    #rolling sum of transaction amounts over last 24 hours
✓ 0s    df["rolling_sum"]= (
            df.rolling(window=86400,on='Time')["Amount"].sum().fillna(0))
```

These features allow even relatively simple models to detect bursty behaviour, rapid retries, or abnormal transaction clustering patterns that are nearly impossible to infer from static features alone.

## Behavioural Deviations and Normalisation

Another high-impact transformation involves measuring deviation from historical norms rather than absolute values. For instance, comparing the current transaction amount to a user's recent spending baseline is often more informative than the amount itself:

```
[15]    #for 7 days for rolling mean amount
✓ 0s    df['mean_amount_7d'] = (
            df.rolling(window=7 * 24 * 60 * 60, on='Time')['Amount']
            .mean()
            .fillna(0)
        )

        #deviation ratio
        df['deviation_ratio'] = (df['Amount'] / df['mean_amount_7d']+1e-6)
```

In credit risk contexts, analogous constructs—such as changes in utilisation ratio or payment volatility—often dominate static borrower attributes. These deviation-based features encode behavioural change, which is central to both fraud emergencies and credit deterioration.

## Preventing Data Leakage in Feature Construction

One of the most common and damaging mistakes in fraud modelling is inadvertent data leakage. Leakage occurs when features are computed using information that would not be available at prediction time, such as future transactions or post-authorisation outsomes. For example, calculating rolling aggregates without enforcing strict time ordering can silently introduce future data into training. Robust feature engineering requires point-in-time correctness, ensuring that every feature reflects only information available up to that transaction.

Although leakage-free features often produce lower offline metrics initially, they generalise far better in production and prevent catastrophic performance collapse after deployment.

## Feature Stability in a Moving Production Environemnt

Production fraud systems operate in adversarial, non-stationary environments. Customer behaviour evolves, fraud strategies adapt, and data pipelines change. Features derived from fragile correlations or rare events tend to decay quickly, forcing frequent retraining and eroding trust. In contrast, stable features such as transaction velocity, normalised spend behaviour, and long-horizon aggregates remain predictive across time and model versions. This stability is especially critical in regulated credit risk settings, where explainability, auditability, and consistency are mandatory.

Ultimately, the enduring lesson from real-world fraud detection is clear: a simpler model trained on strong, stable, and leakage-free features will consistently outperform a complex neural network trained on weak or noisy inputs. Feature engineering is not a preliminary step; it is the core modelling strategy. In high-risk, high-impact domains like fraud and credit risk, simple models plus strong features beat complex models plus weak features every time.

# THE REAL BOTTLENECK IS PRODUCTION, NOT MODELING

Most Machine Learning failures do not originate in notebooks; they emerge quietly after deployment. In offline experiments using datasets like the attached credit card frauds, models often achieve impressive AUC or precision-recall scores using XGBoost, Random Forests, or neural networks trained on Amount, Time, and V1-V28. However, once deployed into a live transaction stream, these same models are exposed to realities that noteooks abstract away: delayed data, schema changes, missing values, and evolving behaviour. The result is a familiar pattern: strong validation metrics during development, followed by gradual performance erosion in production that goes unnoticed until fraud increases or false positives spike.

One of the most common production failure modes is silent data drift. In fraud systems, the statistical properties of core variables, such as Amount or engineered velocity features, can change rapidly due to seasonality, promotions, macroeconomic shifts, or new fraud strategies. Even anonymised PCA features are not immune; their distributions depend on upstream transformations that may evolve over time. Without continuous monitoring, a model may still produce predictions while operating far outside the data regime; it may still produce predictions while operating far outside the regime it was trained on. Feature distribution shifts are especially dangerous because they do not trigger system errors; the pipeline remains "Green" while model quality degrades invisibly.

Another critical bottleneck lies in the data pipelines themselves. In production fraud detection, models are only as reliable as the upstream systems that feed them. Broken joins, delayed ingestion, partial transaction histories, or changes in feature computation logic can all invalidate model assumptions. For example, a rolling feature such as txn_count_1h or amount_sum_24h derived from the Time and Amount fields depends on complete and correctly ordered transaction data. If late-arriving events or pipeline failures truncate windows, the model's perspective, nothing is "wrong"; from the business perspective, trust in the system erodes rapidly.

This is why production machine learning must be treated as a software engineering discipline, not a modelling exercise. Successful fraud systems require robust, versioned data pipelines. consistent feature definitions shared between training and inference; strict latency guarantees to support real-time decisioning; and automated retraining workflows that account for concept drift. Just as importantly, they require monitoring dashboards that track feature distributions, prediction stability, and business KPIs—not just model accuracy. Without these safeguards, models decay quietly in the background until business users stop trusting them, often without anyone being able to pinpoint when or why the failure occurred. In high-stakes domains like fraud and credit risk, operational rigor not algorithmic sophistication, is the true determinant of long-term model success.

# WHERE GENAI FITS AND WHERE IT DOESNOT?

Generative AI has significantly expanded what is possible in modern risk and decisioning systems, particularly when layered on top of traditional machine learning pipelines. In fraud detection and credit risk contexts, transformer-based models enable dense embeddings that capture behavioural similarities across customers, merchants, and transactions. These embeddings can be used for personalisation, anomaly detection, and semantic clustering, surfacing relationships that are difficult to encode manually. Similarly, semantic search over transaction histories, case notes, and customer communications enables investigators and analysts to quickly retrieve relevant context, thereby reducing time-to-decision in fraud operations and credit reviews.

GenAI also unlocks more intuitive interaction paradigms through natural language interfaces. Risk analysts can query transaction patterns, summarise customer behaviour, or explore edge cases using conversational prompts rather than writing SQL or Python. In operational settings, large language models can assist with alert triage, explain model outputs in plain language, and generate structured summaries for downstream workflows. These capabilities meaningfully improve productivity, but they don't change the underlying mechanics of risk prediction itself. The core decision logic, approve, decline, flag, or review, still depends on well-calibrated probabilistic models trained on structured signals.

Crucially, Generative AI doesnot replace the fundamentals of model evaluation, governance, or lifecycle management. Prompt engineering without systematic evaluation quickly becomes fragile and untestable. Outputs can vary with small changes in phrasing, context length, or model version, making reproducibility and auditability

difficult, particularly in regulated environments. Without clearly defined acceptance criteria, monitoring, and rollback strategies, GenAI components introduce opaque risk into systems that demand consistency and traceability. In practice, unmanaged prompts become a new form of technical debt, harder to detect and remediate than traditional model errors.

GenAI systems are most effective when tightly coupled with traditional machine learning signals and structured data. In fraud detection, for example, LLM-generated insights are far more reliable when grounded in features such as transaction velocity, rolling aggregates, historical fraud rates, and expected loss estimates. In credit risk, narrative explanations or customer summaries generated by GenAI should be constrained by model scores, policy rules, and regulatory thresholds. This hybrid approach allows GenAI to enhance interpretation and interaction, while core risk decisions remain anchored in statistically validated models.

Clear business constraints and feedback loops are essential to making GenAI usable at scale. Outputs must be elevated not only for linguistic quality, but for decision impact, bias, and consistency over time. Human-in-the-loop workflows where GenAI suggestions are reviewed, corrected, and fed back into the system are critical for continuous experimentation but undermine trust, particularly in high-stakes domains where errors carry financial, legal, or reputational consequences.

Ultimately, Generative AI is best understood as an accelerator rather than a shortcut. It amplifies the value of strong data foundations, robust ML pipelines, and disciplined governane but it cannot compensate for their absence. Organisations that succeed with GenAI in fraud detection and credit risk will be those that integrate it thoughtfully, treating it as an augmentation layer on top of proven decision systems, rather than a replacement for the fundamentals that make those systems reliable.

# WHAT ACTUALLY WORKS IN THE REAL WORLD?

Across successful machine learning deployments, particularly in fraud detection and credit risk, one pattern consistently stands out: teams start with a business decision, not a model. Effective systems are designed around clear operational questions such as "Should this transaction be approved, challenged, or declined?" or "Should this applicant be offered credit, reviewed manually, or rejected?" High-performing teams explicitly define who consumes the model output, at what point in the workflow it is used, and what concentrated action it triggers. When these elements are unclear, even highly accurate models struggle to gain adoption because their outputs do not map clearly to real decisions.

 Another defining trait of successful implementations is the deliberate use of simple, explainable baselines. Logistic regression scorecards, and shallow tree-based models often form the foundation of production systems because they are the most often form the most powerful, but also because they are interpretable, debuggable, and trusted. In regulated environments such as credit risk, these models make it easier to trust. In regulated environments such as credit risk, these models make it easier to validate assumptions, explain decisions to stakeholders, and meet compliance requirements. More complex models, including gradient boosting or neural networks, are introduced incrementally and only when they demonstrably outperform baselines in ways that matter to the business.

Treating machine learning as a living system rather than a one-time deployment is another critical success factor. Fraud patterns evolve, customer behaviour shifts, and data pipelines change over time, causing even strong models to decay. Teams that succeed build continuous monitoring into their workflows, tracking feature distributions, prediction stability, and business KPIs- and retrain models on a regular cadence. Retraining is not

reactive or ad hoc; it is automated, versioned, and tested, ensuring that updates improve outcomes rather than introduce new risk.

Equally important is the understanding that production ML is inherently cross-functional. Data Scientists may design models, but engineers ensure reliability and latency, domain experts validate assumptions, and business stakeholders define success metrics. In fraud detection, investigators provide critical feedback on false positives and emerging attack patterns; in credit risk, policy teams shape constraints around fairness, explainability, and regulatory compliance. When these groups operate in silos, models fail to reflect real-world constraints and quickly lose relevance.

Ultimately, what works in practice is not cutting-edge algorithms, but disciplined execution. Teams that align models with decisions, start simple, monitor relentlessly, and collaborate across functions consistently outperform those chasing technical novelty. In high-stakes domains like fraud detection and credit risk, sustainable success comes from operational rigour and shared ownership, not from the complexity of the model itself.

# CONCLUSION: MACHINE LEARNING IS A DISCIPLINE, NOT A BREAKTHROUGH

Machine Learning success is rarely the result of a single breakthrough model or a cutting-edge architecture. In real-world deployments, particularly in fraud detection, credit risk, and other high-stakes domains, outcomes are shaped far more by alignment and discipline than by technical novelty. Models only create value when they are tightly aligned with business decisions, grounded in reliable data, and deployed within systems designed to handle change. Without this foundation, even the most sophisticated algorithms quickly become fragile experiments rather than durable solutions.

What separates successful teams is their commitment to execution. This means investing in robust data pipelines, explainable baselines, continuous monitoring, and retraining workflows that reflect how the real world evolves. It also means treating trust as a first-class requirement ensuring that model outputs are understandable, auditable, and consistently reliable for the people who depend on them. In practice, trust is earned not through accuracy metrics alone, but through predictable behaviour over time and clear accountability when systems fail.

The future of machine learning belongs to teams that can turn predictions into decisions reliably, responsibly, and at scale. As ML becomes embedded deeper into operational systems, competitive advantage will come from those who approach it as a disciplined engineering and organisational capability, not a one-time innovation. The winners will not be the teams with the most complex models, but those who can deploy, govern, and sustain ML systems that the business truly relies on.