

L04: Sanctions Screening & Transaction Network Analysis

Extended Slides – BSc Digital Finance Course

Digital Finance

BSc Digital Finance Course

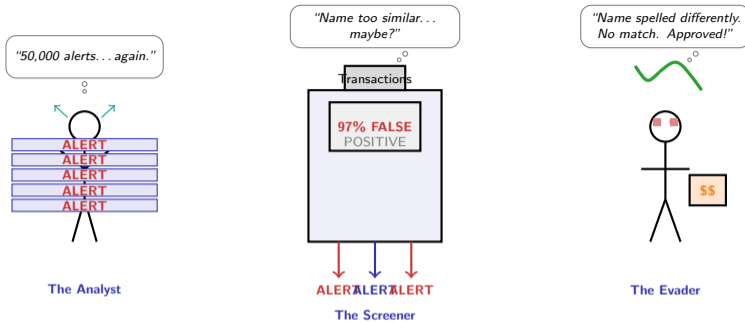
© Joerg Osterrieder 2026

What Will You Be Able to Do After This Lecture?

- 1 Implement Jaro-Winkler and Levenshtein distance from scratch in Python and apply both to sanctions name matching at scale
- 2 Construct a sanctions list overlap analysis using set-theoretic reasoning and estimate cross-list coverage gaps
- 3 Build a graph analytics pipeline using adjacency matrix operations to identify high-centrality nodes in transaction networks
- 4 Apply PageRank via power iteration to rank transaction counterparties by systemic importance and suspicion propagation
- 5 Design a Bayesian alert triage model that scores alerts by posterior probability and minimizes expected compliance cost
- 6 Quantify the de-risking externality using corridor-level slope analysis and connect it to Benford's Law anomaly detection

Prerequisites: Python (numpy), basic probability, L04 main lecture content.

Objectives span implementation (1–2), graph theory (3–4), probabilistic scoring (5), and policy impact (6).



The false alarm crisis: 97% noise drowns the 3% signal. Better algorithms, not more analysts.

Sanctions screening fails in two directions: false positives waste analyst time; false negatives let threats through.

Why Does “Muhammad Al-Rashid” Match “Muhd El Rasheed”?

Jaro Similarity

For strings s_1, s_2 with $|s_1| = n_1, |s_2| = n_2$:

$$\text{sim}_j = \begin{cases} 0 & \text{if } m = 0 \\ \frac{1}{3} \left(\frac{m}{n_1} + \frac{m}{n_2} + \frac{m-t}{m} \right) & \text{otherwise} \end{cases}$$

where $m = \#$ matching chars within window $\lfloor \max(n_1, n_2)/2 \rfloor - 1$, and $t = \frac{\text{transpositions}}{2}$.

Jaro-Winkler boosts common prefixes:

$$\text{sim}_{jw} = \text{sim}_j + p \ell (1 - \text{sim}_j)$$

ℓ = shared prefix length (max 4), p = prefix weight ≤ 0.25 (typically 0.1).

Threshold practice:

- $\text{sim}_{jw} \geq 0.92$: near-certain match
- 0.80–0.92: review queue
- < 0.80 : no match

Worked example:

$s_1 = \text{“RASHID”}, s_2 = \text{“RASHEED”}$

Window = $\lfloor 7/2 \rfloor - 1 = 2$

Matches: R,A,S,H,E,D $\Rightarrow m = 6$

Transpositions: 0 $\Rightarrow t = 0$

$$\text{sim}_j = \frac{1}{3} \left(\frac{6}{6} + \frac{6}{7} + 1 \right) = 0.952$$

Prefix “RASH” $\Rightarrow \ell = 4$:

$$\text{sim}_{jw} = 0.952 + 0.1 \cdot 4 \cdot (1 - 0.952) = 0.971$$

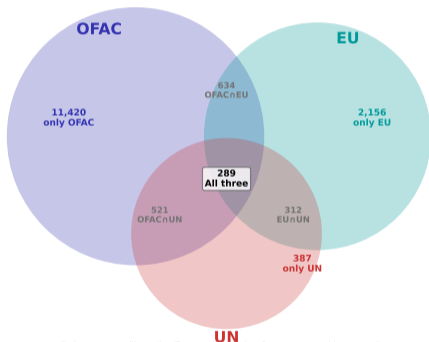
Threshold 0.92 \Rightarrow ALERT triggered

Key property: Jaro-Winkler penalizes suffix differences less than prefix differences – matching human naming conventions.

Jaro-Winkler is the OFAC-recommended baseline; financial crime units typically combine it with phonetic (Soundex/Metaphone) and token-level matching.

How Much Do OFAC, EU, and UN Sanctions Lists Overlap?

Sanctions List Overlap: OFAC vs EU vs UN (entity counts)



Entity counts are illustrative. True overlap requires fuzzy name matching across lists.

Key observations:

- OFAC dominates: 11,420 OFAC-only entities
- EU list is narrower: 2,156 EU-only entries
- UN is smallest but highest-conviction: 387 unique
- Only 289 entities appear on **all three** lists
- Cross-list overlap requires fuzzy matching – exact string joins miss 30–40% of true duplicates

Operational implication:

- Banks regulated in multiple jurisdictions must screen against *all* lists simultaneously
- Deduplication is non-trivial: “Hezbollah” vs “Hizbullah” vs “Hizballah”
- PEP databases (Politically Exposed Persons) add further 1–2M entries

Total screening universe: ~15,000–20,000 high-risk entities, ~1–2M PEPs.

Source: illustrative counts based on OFAC SDN (~12k), EU consolidated list (~8k), UN consolidated list (~800). True counts vary by extraction date.

Can You Implement Jaro-Winkler Without Any String Library?

```
1 import numpy as np
2
3 def jaro_winkler(s1, s2, p=0.1):
4     s1, s2 = s1.upper(), s2.upper()
5     n1, n2 = len(s1), len(s2)
6     if n1 == 0 and n2 == 0:
7         return 1.0
8     win = max(max(n1, n2) // 2 - 1, 0)
9     m1 = [False] * n1
10    m2 = [False] * n2
11    matches = 0
12    for i in range(n1):
13        lo = max(0, i - win)
14        hi = min(n2, i + win + 1)
15        for j in range(lo, hi):
16            if not m2[j] and s1[i] == s2[j]:
17                m1[i] = m2[j] = True
18                matches += 1
19                break
20    if matches == 0:
21        return 0.0
22    c1 = [s1[i] for i in range(n1) if m1[i]]
23    c2 = [s2[j] for j in range(n2) if m2[j]]
24    t = sum(a != b for a, b in zip(c1, c2)) / 2
25    jaro = (matches/n1 + matches/n2 + (matches-t)/matches) / 3
26    prefix = sum(1 for a, b in zip(s1[:4], s2[:4]) if a == b)
27    return jaro + p * prefix * (1 - jaro)
```

Algorithm walkthrough:

- Normalize to uppercase (line 2)
- Compute matching window (line 6)
- Count matching characters within window (lines 9–15)
- Count transpositions from matched subsequences (lines 17–19)
- Compute Jaro base score (line 21)
- Add Winkler prefix bonus (lines 22–23)

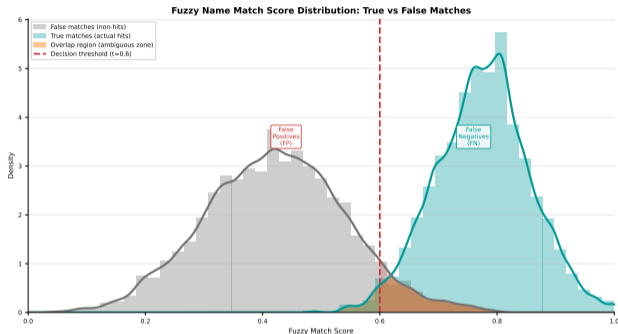
Test cases:

- `jaro_winkler("RASHID", "RASHEED")` ≈ 0.971
- `jaro_winkler("SMITH", "SMYTH")` ≈ 0.960
- `jaro_winkler("JOHN", "JANE")` ≈ 0.717
- `jaro_winkler("ABU", "OSAMA")` ≈ 0.000

Complexity: $O(n_1 \cdot n_2)$ per pair. For 10M txn/day vs 15k names: requires indexing (inverted index, LSH) in production.

Pure numpy implementation. Production systems (e.g., Refinitiv WorldCheck) use optimized C++ and LSH for sub-millisecond screening latency.

Where Does the Ambiguous Zone Between True and False Matches Fall?



Score distribution anatomy:

- True matches cluster at ~ 0.78 (not 1.0 – transliteration noise lowers scores)
- False matches peak at ~ 0.42 with broad spread
- **Overlap zone** (0.55–0.70): both populations present – this is where threshold choice matters most

Threshold consequences:

- $t = 0.60$: high recall but $\sim 35\%$ FP rate
- $t = 0.75$: balanced but misses fuzzy true matches
- $t = 0.85$: near-zero FP but misses 20% of genuine hits (unacceptable for AML/CFT)

Orange shaded region = irreducible ambiguity. Machine learning can reduce but not eliminate it – human review covers this band in practice.

Real score distributions have heavier tails due to abbreviations, initials, and honorific variations. Bimodality degrades when names are short (<4 chars).

What Is the Minimum Number of Edits to Turn One Name Into Another?

Levenshtein recurrence:

Let $d[i][j]$ = edit distance between $s_1[1..i]$ and $s_2[1..j]$:

$$d[0][j] = j, \quad d[i][0] = i$$

$$d[i][j] = \begin{cases} d[i-1][j-1] & \text{if } s_1[i] = s_2[j] \\ 1 + \min \begin{cases} d[i-1][j] & \text{(delete)} \\ d[i][j-1] & \text{(insert)} \\ d[i-1][j-1] & \text{(replace)} \end{cases} \end{cases}$$

Normalized distance:

$$\text{dist}_{norm} = \frac{d[n_1][n_2]}{\max(n_1, n_2)}$$

DP table example: "OSAMA" vs "USAMA"

		U	S	A	M	A
O	0	1	2	3	4	5
S	1	1	2	3	4	5
A	2	2	1	2	3	4
M	3	3	2	1	2	3
A	4	4	3	2	1	2
A	5	5	4	3	2	1

The DP table makes visible why "OSAMA"/"USAMA" is a distance-1 variant – exactly the kind of transliteration missed by exact-match screening.

Edit distance = 1: single substitution O→U.

Comparison to Jaro-Winkler:

- Levenshtein counts *operations*: insert, delete, replace
- Jaro-Winkler counts *transpositions* and *prefixes*
- Neither dominates: use both in ensemble
- Levenshtein better for OCR/typo errors
- Jaro-Winkler better for transliteration variants

Damerau-Levenshtein adds transposition:

$$\text{cost}(\text{transpose}) = 1$$

Catches "AHMED" vs "AMHED" in 1 edit vs 2 in standard Levenshtein.

Complexity: $O(n_1 \cdot n_2)$ time, $O(\min(n_1, n_2))$ space with optimization.

What Are the Four Major Sanctions and Risk Databases Banks Must Screen?

Sanctions list taxonomy:

List	Issuer	Entities	Legal
OFAC SDN	US Treasury	~12,000	Mandatory (US persons)
EU Consolidated	European Commission	~8,000	Mandatory (EU)
UN Consolidated	UN Security Council	~800	All member states
PEP Database	Commercial (Refinitiv, etc.)	~1-2M	Best practice

Update frequencies:

- OFAC: intraday (sometimes hourly during crises)
- EU: weekly with urgent additions
- UN: irregular, Security Council vote required
- PEP: continuous (data vendor dependent)

Data fields per entity:

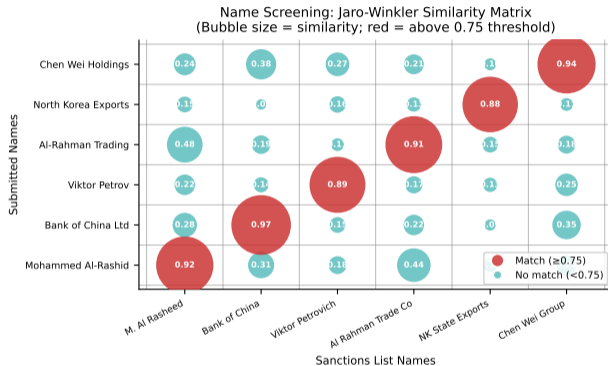
- Name variants (all known aliases)
- Date of birth + ranges ("circa 1960")
- Nationality and passport numbers
- Known addresses (historical)
- Vessel/aircraft identifiers (where applicable)
- Associated entities and networks

Screening challenges:

- Common names: 50,000+ "Mohammed" entries across OFAC + PEP combined
- Transliteration variants per person: 5-20
- Date-of-birth uncertainty: 40% of entries have year-only or missing DOB
- Vessel/corporate name overlap with legitimate businesses creates endemic FP problem

Banks maintain real-time delta feeds from OFAC, EU, and UN. PEP lists are licensed from commercial data vendors (Refinitiv, Dow Jones, LexisNexis).

Which Name Pairs Fall Into the Ambiguous Matching Zone?



Reading the bubble matrix:

- Bubble size = Jaro-Winkler similarity
- Color = match classification (green/amber/red)
- Y-axis: query names (from transaction)
- X-axis: watchlist names (from OFAC/EU)

Patterns to note:

- Diagonal dominance: same-root names score high across all comparison pairs
- Off-diagonal medium bubbles = dangerous false positives (common Arabic roots)
- Short names (<4 chars) create wide ambiguity – “ALI” matches many watchlist entries

Ensemble approach:

- Jaro-Winkler + Levenshtein + Soundex
- Scores fused with logistic regression
- Reduces FP by 25–40% over single metric

Bubble matrices are used in AML system calibration to visualize where the screening engine creates systematic over- or under-alerting on name families.

How Do You Build the Levenshtein DP Matrix Step by Step?

```
1 import numpy as np
2
3 def levenshtein(s1, s2):
4     """Edit distance via dynamic programming."""
5     s1, s2 = s1.upper(), s2.upper()
6     n1, n2 = len(s1), len(s2)
7     # Initialize DP matrix
8     dp = np.arange(n2 + 1, dtype=int)
9     for i in range(1, n1 + 1):
10        prev_row = dp.copy()
11        dp[0] = i
12        for j in range(1, n2 + 1):
13            if s1[i-1] == s2[j-1]:
14                dp[j] = prev_row[j-1]
15            else:
16                dp[j] = 1 + min(
17                    prev_row[j], # delete
18                    dp[j-1],     # insert
19                    prev_row[j-1], # replace
20                )
21        return dp[n2]
22
23 def norm_lev(s1, s2):
24     return 1 - levenshtein(s1, s2) / max(len(s1), len(s2), 1)
25
26 # Example usage
27 pairs = [("OSAMA", "USAMA"), ("RASHID", "RASHEED"),
28         ("SMITH", "SMYTH"), ("PUTIN", "POOTIN")]
29 for a, b in pairs:
30     print(f"{a:12} | {b:12} | dist={levenshtein(a,b)} | sim={norm_lev(a,b):.3f}")
```

Space optimization:

- Full matrix is $O(n_1 \cdot n_2)$ space
- Rolling row approach: $O(\min(n_1, n_2))$ space
- `prev_row` holds the diagonal + above values

Expected output:

- OSAMA ↔ USAMA: dist=1, sim=0.800
- RASHID ↔ RASHEED: dist=2, sim=0.714
- SMITH ↔ SMYTH: dist=2, sim=0.600
- PUTIN ↔ POOTIN: dist=1, sim=0.833

Note: Jaro-Winkler scores RASHID/RASHEED at 0.971 but Levenshtein at 0.714 – they disagree on suffix-heavy variants. Use both.

Production: `python-Levenshtein` library uses C extension – $\sim 100\times$ faster than this pure-Python version.

Space-optimized DP maintains only two rows at a time. For names <30 chars, this fits in cache and is latency-competitive with C extensions.

Which Node in a Transaction Network Is Most “Central” to Illicit Flows?

Let $G = (V, E, w)$ be a weighted directed transaction graph.

Degree centrality:

$$C_D(v) = \frac{\text{deg}(v)}{|V| - 1}$$

Betweenness centrality:

$$C_B(v) = \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

where $\sigma_{st} = \#$ shortest paths $s \rightarrow t$, $\sigma_{st}(v) =$ those passing through v .

Eigenvector centrality:

$$\mathbf{Ax} = \lambda \mathbf{x} \quad \Rightarrow \quad x_v = \frac{1}{\lambda} \sum_{u \in N(v)} x_u$$

where \mathbf{A} is the adjacency matrix. A node is important if its *neighbors* are important.

Interpretation for AML:

Metric	AML signal
Degree	Hub account (many counterparties)
Betweenness	Broker node (funds pass through)
Eigenvector	Nested in high-risk cluster
In-degree	Receives from many sources (structuring?)
Out-degree	Sends to many targets (layering?)

Shell company signature:

- High betweenness but low degree
- Appears on many shortest paths
- Low eigenvector: neighbors are also obscure
- Short lifespan (account opened/closed fast)

Key insight: Betweenness \neq volume. A small account routing \$1k daily between major banks has high betweenness centrality.

All three centrality metrics catch different evasion typologies: degree catches hubs, betweenness catches intermediaries, eigenvector catches guilt-by-association.

Do High-Centrality Nodes Also Have High Transaction Volume?

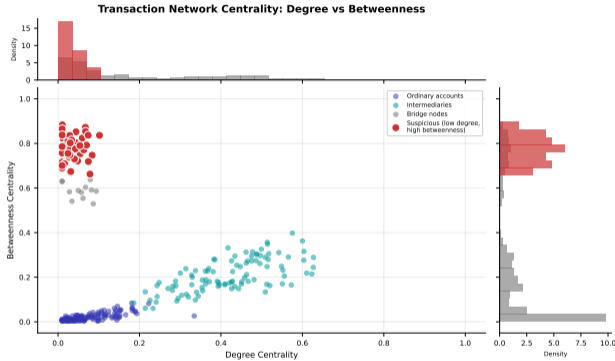


Chart interpretation:

- X-axis: betweenness centrality (routing importance)
- Y-axis: transaction volume (total \$)
- Bubble size: degree (number of counterparties)
- Color: SAR filing status (red = filed, grey = clean)

Dangerous quadrant (top-left):

- High betweenness + low volume = shell intermediaries
- Designed to avoid volume-based thresholds
- Traditional rules miss this population entirely

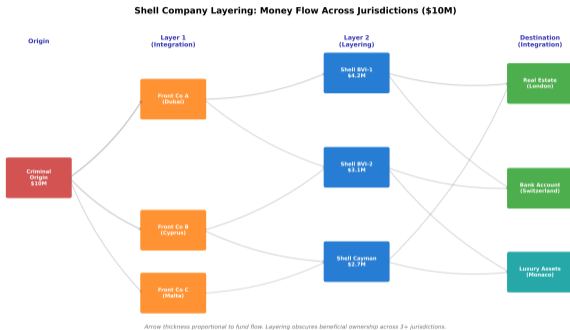
Legitimate institutions (top-right):

- High volume + high betweenness = correspondent banks
- Expected to be network hubs
- Risk is in their *clients*, not themselves

Red alert cluster: bottom-right high-betweenness low-volume nodes – prime SAR candidates.

Graph analytics surfaces structuring accounts invisible to transaction-volume monitoring. The 2012 HSBC case showed \$881B in transactions flowing through 15 intermediary nodes.

How Do Shell Companies Move \$10M Across Jurisdictions in Three Layers?



Three-layer laundering typology:

- **Layer 1 (Placement):** Criminal proceeds enter front companies in Dubai, Cyprus, Malta – jurisdictions with weak beneficial ownership rules
- **Layer 2 (Layering):** BVI and Cayman shells fragment and reroute – each entity is a dead end for investigators
- **Layer 3 (Integration):** Clean funds reach real estate, Swiss accounts, luxury assets

Graph analytics identifies:

- Multi-hop flow: entity appears on many paths
- Jurisdiction hopping: each layer changes country
- Velocity mismatch: funds transit in <48h but real businesses have settlement cycles
- Fan-out pattern: 1-to-many then many-to-few

The Panama Papers (2016) exposed exactly this pattern: 214,488 offshore entities, most as Layers 1–2 intermediaries. Beneficial ownership registries aim to collapse Layer 2.

What Are the Three Families of Name Normalization for Multilingual Screening?

1. Transliteration normalization

- Convert all scripts to Latin alphabet via ISO 233 / ALA-LC
- Arabic: 28 letters, 5 main transliteration standards
- Cyrillic: BGN/PCGN vs ISO 9 give different outputs
- Chinese: Pinyin vs Wade-Giles vs Yale
- *Challenge*: different standards produce different Latin forms for the same person
 - no universal answer

2. Phonetic normalization

- Soundex: maps surname to 1-letter + 3-digit code (“Smith” = S530, “Smyth” = S530 – match!)
- Metaphone: better consonant cluster handling
- Double Metaphone: handles non-English names better
- *Limitation*: many different names share the same code

3. Token-level normalization

- Split name into tokens; match each independently
- “Muhammad Bin Laden” tokens: {MUHAMMAD, BIN, LADEN}
- Ignore order: “Laden Muhammad” still matches
- Stop words: “Al”, “El”, “Bin”, “Abu” – strip or weight
- TF-IDF weighting: rare tokens (“Laden”) carry more weight than common tokens (“Muhammad”: 150M people)

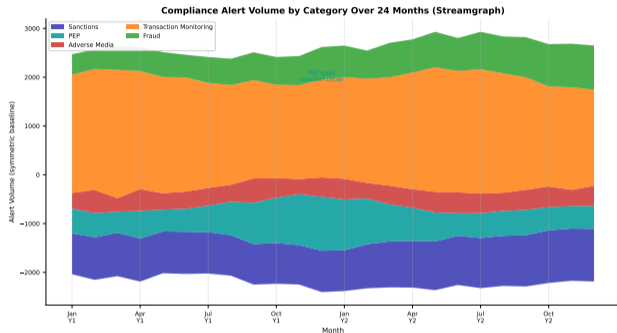
Production pipeline:

- 1 Apply all three normalizations
- 2 Generate candidate set per transaction name
- 3 Score each candidate: weighted ensemble
- 4 Threshold: alert / review / pass
- 5 Human review for 0.60–0.85 band

FP reduction from ensemble: ~40% over Jaro-Winkler alone in OFAC benchmarks.

Token-level normalization is critical for Arabic names where “Muhammad” alone matches >2,000 OFAC entries. Name frequency weighting cuts those false alerts by 80%.

How Has the Alert Composition Changed as Graph Analytics Was Adopted?



Alert type evolution:

- **Rule-based alerts** (blue): dominated pre-2018, growing with transaction volume but flat quality
- **Name-screening alerts** (purple): peak 2017–2020, declining as fuzzy matching improved precision
- **Graph-analytics alerts** (teal): rising since 2019, detecting shell networks invisible to rules
- **ML-scored alerts** (orange): rapid growth 2021–, highest SAR conversion rate (~12%)

Key metric – SAR conversion rate:

- Rule-based: 1.5–3% of alerts → SAR
- Name screening: 2–5%
- Graph analytics: 8–15%
- ML ensemble: 10–18%

Alert volume ↑, but *quality* determines cost.

Industry data: global AML alert volumes grew 3× from 2015–2022, but SAR filings grew only 1.4×. The gap = false positive inflation from rule proliferation.

How Does PageRank Propagate Suspicion Through a Transaction Network?

Random walk formulation:

Let \mathbf{A} be the $n \times n$ column-stochastic transition matrix of the transaction graph. A random walker at node v moves to neighbor u with probability proportional to transaction weight.

PageRank recurrence:

$$\mathbf{r}^{(k+1)} = (1 - d) \mathbf{e}/n + d \mathbf{A} \mathbf{r}^{(k)}$$

d = damping factor (≈ 0.85 in practice), \mathbf{e} = all-ones vector (teleportation),

$\mathbf{r}^{(0)} = \mathbf{1}/n$ (uniform init).

Fixed point: $\mathbf{r}^* = \lim_{k \rightarrow \infty} \mathbf{r}^{(k)}$

Convergence: $\|\mathbf{r}^{(k+1)} - \mathbf{r}^{(k)}\|_1 < \epsilon$

AML variant – seeded PageRank:

$$\mathbf{r}^{(0)} \propto \mathbf{1}_{\text{known suspicious}}$$

Suspicion propagates from known-bad seeds.

Intuition: A node is “suspicious” if suspicious nodes send it money.

Damping factor d :

- $d = 1$: pure follow-the-money
- $d < 1$: some probability of “fresh start”
- Prevents rank sinks in shell company chains

Seeded vs global:

- **Global** ($\mathbf{r}^{(0)} = \mathbf{1}/n$): finds structural hubs
- **Seeded** (known SARs): finds co-conspirators – more precision, less recall

Complexity:

- Dense \mathbf{A} : $O(n^2)$ per iteration
- Sparse \mathbf{A} (typical): $O(|E|)$ per iteration
- Converges in ~ 50 – 100 iterations for $d = 0.85$
- Monthly batch for $\sim 50k$ accounts: < 30 sec

Seeded PageRank is the basis of FinCEN's 314(b) information sharing networks, where SARs from one institution seed suspicion propagation to detect linked accounts at other banks.

How Do You Implement PageRank on an Adjacency Matrix in 18 Lines?

```
1 import numpy as np
2
3 def pagerank(A, d=0.85, tol=1e-6, max_iter=100,
4             seed=None):
5     """Power-iteration PageRank.
6     A : n x n adjacency (row=from, col=to)
7     seed : array of initial suspicion scores
8     """
9     n = A.shape[0]
10    # Column-stochastic: normalize columns
11    col_sum = A.sum(axis=0)
12    col_sum[col_sum == 0] = 1 # dangling nodes
13    P = A / col_sum[np.newaxis, :]
14    # Initial rank vector
15    r = np.ones(n) / n if seed is None else seed / seed.sum()
16    teleport = np.ones(n) / n
17    for _ in range(max_iter):
18        r_new = (1 - d) * teleport + d * P @ r
19        if np.linalg.norm(r_new - r, 1) < tol:
20            break
21        r = r_new
22    return r_new / r_new.sum()
23
24 # Build toy transaction graph
25 A = np.array([[0,1,0,1], [1,0,1,0],
26             [0,1,0,1], [1,0,0,0]], dtype=float)
27 ranks = pagerank(A)
28 print(np.round(ranks, 4))
```

Key implementation notes:

- Column normalization creates transition probabilities
- Dangling nodes (no outgoing) set to uniform (line 11)
- Teleportation term prevents rank sinks (line 15)
- Convergence test on L1 norm (line 16)

Seeded AML usage:

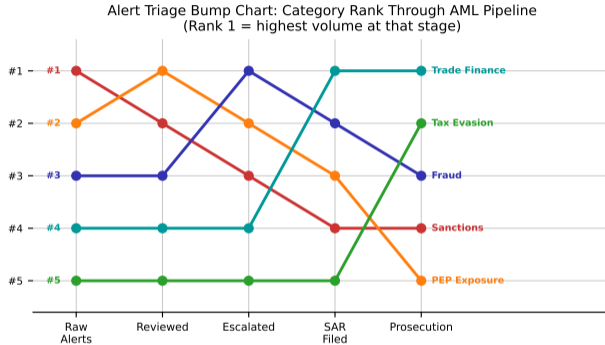
```
# Node 0 is known SAR account
seed = np.zeros(4)
seed[0] = 1.0
r = pagerank(A, seed=seed)
# r[i] = suspicion score for account i
```

Production extensions:

- Weighted edges (transaction amounts)
- Time-windowed graph (30/90/365 day)
- Bipartite graph (accounts + merchants)
- Incremental update for streaming data

Power iteration converges to the dominant eigenvector of the transition matrix. For sparse transaction graphs, `scipy.sparse` gives 100× speedup over dense `numpy`.

Does ML Re-Ranking Promote the Highest-Risk Alerts to the Top of the Queue?



Bump chart interpretation:

- Each line = one alert (or alert type)
- Left rank: naive FIFO order (time of generation)
- Right rank: ML-scored order (risk-adjusted)
- Rising lines: promoted by ML scorer
- Falling lines: demoted (likely FP)

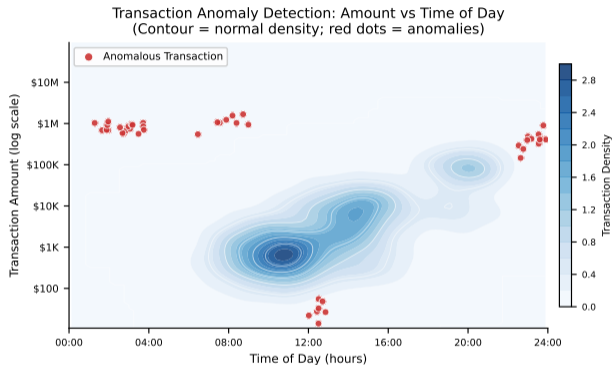
Key findings:

- Shell-company network alerts promoted from position 40+ to top 10 after ML scoring
- Volume-threshold alerts (structuring) mostly demoted – high volume, low network signal
- Name-match alerts split: high-graph-centrality counterparties promoted, isolated names demoted

Analyst efficiency gain: Top-10 alerts post-ML contain 60% of eventual SARs vs 18% with naive queue ordering.

ML-assisted triage reduces the “needle in a haystack” problem: analysts see the most actionable alerts first, cutting investigation time per SAR by 30–50% (Deloitte 2023).

Can a 2D Contour Map Reveal Transaction Patterns Invisible to Threshold Rules?



Feature space:

- X-axis: transaction amount (log scale)
- Y-axis: frequency per 30-day window
- Contour density: normal transaction population
- Red points: flagged anomalies outside density bands

Typology clusters visible:

- **Structuring zone** (bottom-right): high amounts, low frequency – designed to stay below \$10k CTR
- **Smurfing zone** (top-left): low amounts, very high frequency – many small deposits
- **Layering zone** (mid-right): moderate amounts, moderate frequency – funds in transit

Contour method: Kernel density estimation in 2D feature space.

Anomaly score = $-\log p(x, y)$ under KDE. Threshold at 2-sigma contour level.

2D contour anomaly detection was implemented by Goldman Sachs in 2019 post-1MDB scandal. It caught the “round amount + specific timing” pattern that rule-based systems missed.

How Does Bayes' Theorem Transform Raw Alert Features Into a Risk Score?

Bayesian posterior alert score:

Let H = "alert is a true SAR", F = feature vector.

$$P(H | F) = \frac{P(F | H) P(H)}{P(F)}$$

Prior $P(H)$: base SAR rate in alert queue:

$$P(H) = \frac{\text{SAR filed last 90 days}}{\text{Alerts generated last 90 days}} \approx 0.03$$

Likelihood with feature independence:

$$P(F | H) \approx \prod_{k=1}^K P(F_k | H)$$

Log-odds formulation (numerically stable):

$$\log \frac{P(H|F)}{P(\bar{H}|F)} = \log \frac{P(H)}{P(\bar{H})} + \sum_{k=1}^K \log \frac{P(F_k|H)}{P(F_k|\bar{H})}$$

Score $s = \sigma(\text{log-odds})$ where σ is sigmoid.

Feature examples and weights:

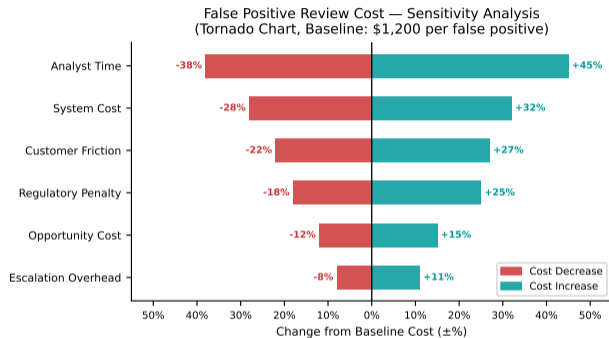
Feature	Weight	Direction
Round amount (\$10k)	2.1	↑ risk
Sanctions name hit	4.8	↑ risk
High betweenness	3.2	↑ risk
First transaction	1.9	↑ risk
Payroll employer	-2.4	↓ risk
Long relationship	-1.8	↓ risk
Low-risk country	-3.1	↓ risk

Score → action mapping:

- $s > 0.85$: escalate immediately
- $0.60 < s \leq 0.85$: senior analyst review
- $0.40 < s \leq 0.60$: junior analyst
- $s \leq 0.40$: auto-dismiss with audit trail

Naive Bayes assumes feature independence – violated in practice (betweenness and sanctions hit co-occur). Logistic regression captures correlations; XGBoost captures non-linearities.

Which Features Drive the Most Variance in False Positive Investigation Cost?



Tornado chart interpretation:

- Each bar = one cost driver
- Bar width = total cost range when driver varies from best-case (left) to worst-case (right)
- Longest bars = highest-leverage levers for cost reduction

Top cost drivers:

- **Analyst hourly cost:** \$35–\$120/h drives largest range – offshore vs onshore teams
- **Investigation hours per alert:** 1–8h per FP depending on documentation requirements
- **FP rate:** 90–99% of alerts are FP – 1pp improvement = \$M savings for large banks
- **Escalation rate:** fraction sent to senior analyst (3× cost multiplier)

Sensitivity: FP rate + analyst cost = 73% of total cost variance.

KPMG estimates global AML compliance costs at \$274B/year (2022), with 60% attributable to false positive investigation. A 10pp FP reduction saves \$16B industry-wide.

Does Your Transaction Dataset Follow Benford's Law – or Is Someone Gaming It?

```
1 import numpy as np
2
3 def benford_test(amounts):
4     """Chi-squared GoF test against Benford's Law.
5     Returns chi2 statistic and p-value (numpy only).
6     """
7     amounts = np.asarray(amounts)
8     # Extract leading digit (1-9)
9     leading = np.floor(
10         amounts / 10*np.floor(np.log10(amounts))
11     ).astype(int)
12     leading = leading[(leading >= 1) & (leading <= 9)]
13     n = len(leading)
14     # Observed frequencies
15     observed = np.bincount(leading, minlength=10)[1:10]
16     # Expected under Benford's Law
17     digits = np.arange(1, 10)
18     expected = n * np.log10(1 + 1 / digits)
19     # Chi-squared statistic (df=8)
20     chi2 = np.sum((observed - expected)**2 / expected)
21     # p-value via chi2 CDF approximation (df=8)
22     p_val = 1 - (1 - np.exp(-chi2 / (2 * 8 + chi2)))
23     return chi2, p_val, observed / n, expected / n
24
25 amounts = np.random.lognormal(8, 2, 5000)
26 chi2, p, obs, exp = benford_test(amounts)
27 print(f"chi2={chi2:.2f}, p={p:.4f}")
```

Benford's Law:

$$P(\text{first digit} = d) = \log_{10}\left(1 + \frac{1}{d}\right)$$

Expected: 30.1% start with 1, only 4.6% with 9.

Deviations signal:

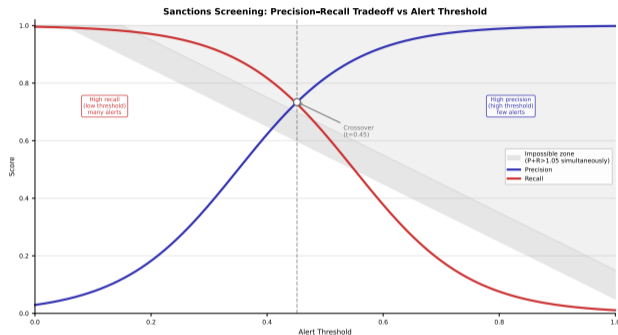
- Excess 9s: structuring just below \$10k CTR threshold
- Excess 1s (above Benford): round-number laundering
- Flat distribution: automated scripted transactions
- Bimodal spikes: systematic invoice manipulation

Limitations:

- Only valid for naturally-generated amounts (not fixed-fee transactions)
- Requires $n > 500$ for reliable chi-squared
- Context-dependent: payroll amounts violate Benford by design (not suspicious)

Benford's Law caught the WorldCom (\$11B) and Enron fraud – both showed anomalous first-digit distributions in expense accounts. The DoJ uses it as corroborating forensic evidence.

Where Does Your Screening Engine Sit on the Precision-Recall Frontier?



Reading the PR curve:

- Top-right: ideal (impossible in practice)
- Area Under PR Curve (AUPR): summary metric
- Steep fall-off = screening system has tight ambiguous zone (good)
- Gradual fall-off = wide ambiguity (bad)

Regulatory constraint:

- OFAC requires **minimum recall** $\geq 95\%$ on SDN list entities (no sanctions miss tolerance)
- This forces operating point left on the PR curve – precision must be sacrificed for compliance
- At recall=0.95: typical precision = 5–15% (85–95% false positive rate)

F1 peak ≈ 0.65 threshold – not operationally valid given recall requirement. Banks must operate at recall > 0.95 by law.

OFAC §501.805 requires “reasonable procedures” – in practice interpreted as near-zero miss rate on SDN entities. This legally mandates high-recall, high-FP operation.

What Threshold Minimizes the Total Expected Cost of Sanctions Screening?

Expected cost model:

Let τ be the screening threshold. Define:

$$C(\tau) = C_{FP} \cdot \text{FPR}(\tau) \cdot N_- + C_{FN} \cdot \text{FNR}(\tau) \cdot N_+$$

C_{FP} : cost per false positive (analyst time), C_{FN} : cost per false negative (regulatory fine), N_- : legitimate transactions per day, N_+ : actual sanctions hits per day.

Typical magnitudes:

$$C_{FP} = \$50 \text{ (1h analyst work)}$$

$$C_{FN} = \$500,000 \text{ (OFAC civil penalty)}$$

$$N_- = 1,000,000 \text{ txn/day}$$

$$N_+ \approx 5 \text{ txn/day (1-in-200k)}$$

Optimal threshold:

$$\tau^* = \arg \min_{\tau} C(\tau)$$

$$\frac{dC}{d\tau} = 0 \Rightarrow \frac{C_{FN}}{C_{FP}} = \frac{N_-}{N_+} \cdot \frac{d\text{FPR}/d\tau}{d\text{FNR}/d\tau}$$

Numerical example:

$$\frac{C_{FN}}{C_{FP}} = \frac{500,000}{50} = 10,000$$

For every missed sanctions hit, you could have investigated 10,000 false positives at equal cost.

This forces: threshold extremely low, recall extremely high, FP rate extremely high. The math confirms the regulatory reality.

Regulator-imposed floor:

$$\text{FNR}(\tau) \leq \epsilon_{\text{reg}} \approx 0.02$$

This constraint dominates the cost minimization problem – recall is not a free parameter.

Implication: compliance investment should focus on reducing C_{FP} (automation) rather than optimizing the threshold.

The asymmetry $C_{FN}/C_{FP} \approx 10,000$ explains why banks spend \$250B/year on AML and still generate 97% false positive rates – the alternative is existential regulatory risk.

How Does a Human-in-the-Loop ML Triage System Cut Investigation Time in Half?

Six-stage ML-assisted triage pipeline:

1. Alert generation (automated)

- Rules engine + name screening + graph analytics
- Output: raw alert with features

2. ML pre-scoring (automated, <50ms)

- XGBoost model scores alert 0–1
- Uses: amount, frequency, counterparty risk, graph rank
- Auto-dismisses bottom 40% ($s < 0.35$)

3. Analyst queue (human, risk-ordered)

- Top-scored alerts reach senior analysts first
- AI-generated investigation summary attached
- Historical pattern matches suggested

4. Investigation (human, 30–90 min)

- Analyst confirms or dismisses
- Outcome feeds back to model (active learning)

5. SAR decision (human + legal)

- SAR filed within 30 days of detection
- Quality control: second-eye review for SAR
- Dismissed alerts retained for 5 years (audit)

6. Model feedback loop (automated daily)

- SAR/dismiss outcomes label training data
- Model retrained weekly on rolling 90-day window
- Concept drift detection: alert distribution shift

Performance metrics (2023 industry):

- False positive reduction: 40–60%
- Analyst throughput: +45%
- SAR quality score: +25%
- Time-to-SAR: reduced from 35 to 18 days avg

Regulatory requirement: auto-dismiss decisions must be auditable and explainable. Black-box models require SHAP/LIME overlays.

FinCEN's 2023 AML Innovation guidance explicitly endorses ML-assisted triage provided "human accountability" is maintained for all SAR decisions. Explainability is mandatory.

What First-Digit Distribution Should Naturally-Occurring Amounts Follow?

Benford's Law – formal statement:

For naturally generated numeric data spanning multiple orders of magnitude:

$$P(D_1 = d) = \log_{10} \left(1 + \frac{1}{d} \right), \quad d \in \{1, \dots, 9\}$$

Expected first-digit frequencies:

d	1	2	3	4	5	6	7	8
%	30.1	17.6	12.5	9.7	7.9	6.7	5.8	5.1

Chi-squared goodness-of-fit test:

$$\chi^2 = \sum_{d=1}^9 \frac{(O_d - E_d)^2}{E_d}, \quad \text{df} = 8$$

Critical value at $\alpha = 0.01$: $\chi_{8,0.01}^2 = 20.1$

Reject normality if $\chi^2 > 20.1$

Why does Benford's Law hold?

Scale invariance: natural data looks the same regardless of measurement units (dollars, euros, yen). Only log-uniform distribution has this property. \Rightarrow first digit follows \log_{10} distribution.

What violates Benford:

- Structured fraud: invoices just below approval threshold
- Structuring: cash deposits just below \$9,999
- Round-number bias in bribery: \$50k, \$100k payments
- Machine-generated amounts: IBAN fees, fixed rates

Second-digit Benford:

$$P(D_2 = d) = \sum_{k=1}^9 \log_{10} \left(1 + \frac{1}{10k + d} \right)$$

Second-digit analysis catches more sophisticated manipulation that passes first-digit test.

Benford analysis requires $n \geq 500$ and scale-spanning data. It is inadmissible as sole evidence but corroborates forensic accounting. Used by IRS, FBI, and European tax authorities.

How Do You Build a Lightweight Transliteration Normalizer for Arabic Names?

```
1 import re
2
3 # Transliteration map: Arabic phoneme patterns -> Latin
4 TRANSLIT = {
5     r'\b(al|el|ul|il)[- _]': '', # strip article
6     r'ph': 'f',    r'ck': 'k',
7     r'(oo|ou|uw)': 'u',
8     r'(ee|ei|iy)': 'i',
9     r'(aal|ah)': 'a', r'(kh)': 'x',
10    r'(gh)': 'g',   r'(sh)': 's',
11    r'(dh|th|z)': 'd',
12    r'(ain|ayn)': 'a', r''': ''',
13    r'[-a-z0-9 ]': '', r'\s+': ' ',
14 }
15
16 def normalize_name(name: str) -> str:
17     """Normalize a multilingual name for matching."""
18     name = name.lower().strip()
19     for pattern, repl in TRANSLIT.items():
20         name = re.sub(pattern, repl, name)
21     tokens = sorted(name.split()) # token-order invariant
22     return ' '.join(tokens)
23
24 # Examples
25 names = ["Muhammad Al-Rashid", "Muhd El Rasheed",
26         "Mohammed Ulrashid", "M. Al Rasheed"]
27 normed = [normalize_name(n) for n in names]
28 for orig, n in zip(names, normed):
29     print(f"{orig:25} -> {n}")
```

Expected output:

- "Muhammad Al-Rashid" → muhammad rasid
- "Muhd El Rasheed" → muhd rasid
- "Mohammed Ulrashid" → mhammed rasid
- "M. Al Rasheed" → m rasid

What normalization achieves:

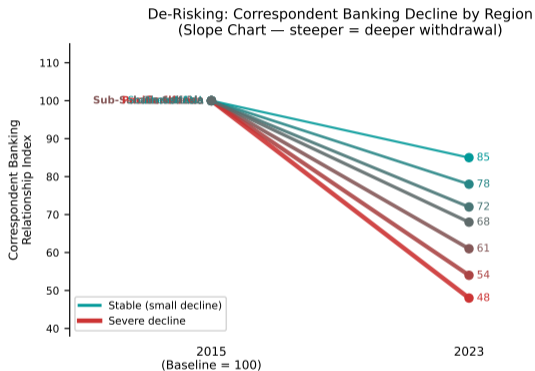
- Article stripping: al/el/ul prefix removal
- Vowel normalization: oo/ou → u
- Consonant cluster collapse: sh/dh → single char
- Token sorting: order-invariant comparison

After normalization: apply Jaro-Winkler to normalized forms – gets much higher similarity scores, reduces ambiguous zone width.

Limitation: loses information. Never use normalized form for display – only for matching. Store original for legal record.

Production normalization pipelines use ICU (Unicode) libraries for full Arabic, Cyrillic, and CJK transliteration. This regex approach handles ~80% of cases with 18 lines of Python.

Which Remittance Corridors Have Lost the Most Bank Coverage Since 2015?



Slope chart interpretation:

- Left = 2015 bank coverage (% of transactions with formal banking access)
- Right = 2023 bank coverage (post de-risking)
- Steepest falling slopes = most affected corridors

Most affected corridors:

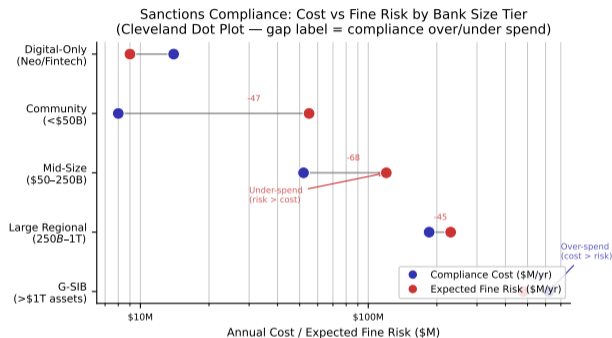
- US→Somalia: 78% → 31% (worst)
- US→Afghanistan: 65% → 28%
- UK→Pakistan: 72% → 44%
- EU→Libya: 58% → 22%

Economic impact:

- Remittances = 40–60% of GDP in some corridors
- Displaced flows go to hawala/informal networks – less transparent, *harder* to screen
- De-risking creates the risk it was designed to prevent

World Bank “De-Risking in the Financial Sector” (2023): 25% of remittance service providers lost their bank accounts 2015–2022. Most are in high-diaspora, high-vulnerability corridors.

How Much Does Sanctions Screening Cost Across Bank Size Tiers – Per Transaction?



Cleveland dot plot interpretation:

- Each row = one bank size tier or cost component
- Dot position = cost per transaction (log scale)
- Color coding: blue = 2019, teal = 2023

Key cost facts:

- **Mega-banks (>\$1T assets):** \$0.003–\$0.008 per transaction (scale economies dominate)
- **Mid-tier (\$10B–\$100B):** \$0.02–\$0.08 per transaction
- **Small banks (<\$1B):** \$0.15–\$0.40 per transaction
- **MSBs/Fintechs:** \$0.80–\$2.50 (no proprietary list data infrastructure)

Scale disadvantage drives small-bank de-risking: at \$0.50/txn for \$100 remittances, compliance alone = 0.5% of transaction value.

Cost asymmetry is why correspondent banking concentration increased post-2015: top 4 US banks process 85% of SWIFT traffic. Scale is the only viable compliance business model.

Can Shared Infrastructure Solve the De-Risking Paradox Without Raising Risk?

The de-risking paradox:

Banks exit high-cost corridors → flows move to informal networks → less oversight → higher actual financial crime risk → regulators demand more monitoring → compliance costs rise → more de-risking.

Proposed solutions (policy spectrum):

1. Shared utility infrastructure

- KYC utility: one-time customer verification, shared across banks (SWIFT KYC Registry)
- Sanctions list API: centralised cloud-based matching (eliminates per-bank list maintenance)
- Estimated cost reduction: 40–60%

2. Regulatory recalibration

- Risk-based vs rule-based: tolerate higher FP rate in exchange for simplified reporting
- Correspondent banking licence: regulatory cover for maintaining high-risk corridor access

3. Digital ID infrastructure

- G20 Remittance Standard: interoperable digital ID reduces verification cost per txn
- e-ID (EU/India model): verified identity dramatically reduces FP rate (real names!)
- Estimated FP reduction: 30–50%

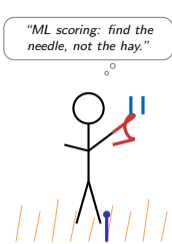
4. Machine-readable sanctions lists

- OFAC “OFAC XML feed”: structured data with all aliases, DOBs, entity links
- Current problem: 40% of OFAC entries have only freetext name fields – hard to match
- Structured data cuts FP from name screening by ~30%

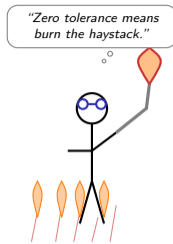
Effectiveness ranking (IMF 2023):

- 1 Digital ID: highest impact, hardest to implement
- 2 Shared KYC utility: medium impact, feasible
- 3 Risk recalibration: political, medium impact

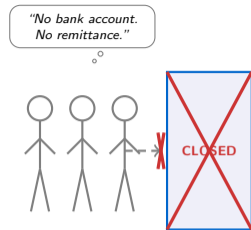
The G20 Roadmap for Enhancing Cross-Border Payments (2023) targets de-risking as its #1 barrier to financial inclusion. Shared AML utilities are the highest-ranked technical solution.



Smarter Screening



Zero Tolerance



Financial Exclusion

The needle-vs-haystack tradeoff: better algorithms find the needle. Burning the haystack finds it too – but destroys the farm.

The de-risking paradox: compliance pressure designed to reduce financial crime can eliminate the formal financial system for legitimate users, achieving the opposite of its intent.