# PAYMENT FRAUD (TRANSACTION LEVEL)

**ONLINE METRICS:**
- ✓ i) Recall/Detection Rate: % of actual fraud caught
- ✓ ii) False +ve rate: % of legitimate transactions incorrectly flagged as fraudulent
- iii) Precision

Trade-off: FPR vs Recall: Lower FPR reduces friction but may miss fraud; Higher recall catches more fraud but increases false +ve

- ✓ iv) Revenue loss
- ✓ v) Manual review queue size
- vi) Flag rate: (until labels are available)

**LABEL CRITERIA:**
- i) Initial labels: chargebacks, refund reason codes
- ii) Due to label delay (45 to 60 days) and evolving "fraud" def"/tactics:
  → Manual review on rule-based & ML model flagged + customer reported fraud
- iii) Feedback loop ISSUE: Transactions blocked by model never get true labels even blocked by model ≠ manual review of a sample

**TRAINING DATA:** Imbalanced dataset

**FEATURES:**
- → Data leakage: ensuring model doesn't learn from post-transaction signals
- → Privacy concerns: PII/sensitive data
- → Feature freshness: real time features valuable but has cost attached (due to evolving fraud tactics) (infra, simple feat engg)
- → Feature drops: i) Velocity features: ratio of purchases in last hr to last 30 days
  (transaction-level so account-level info may not be used) {Trans. history based on credit card, email etc} ratio of amt spent in last hr to last 30 days
  - ii) Affinity features on top of ratio/counter based features: ratio of purchases in last 1 hrs to last 24 hr in a) location b) credit card
  - iii) Reputation features: "email domain" "IP" "phone no." "device id" (emailage, telesign, maxmind, nData)
  - iv) Account data (if it can be linked): account age, gender, age etc; account level transaction history + Profile
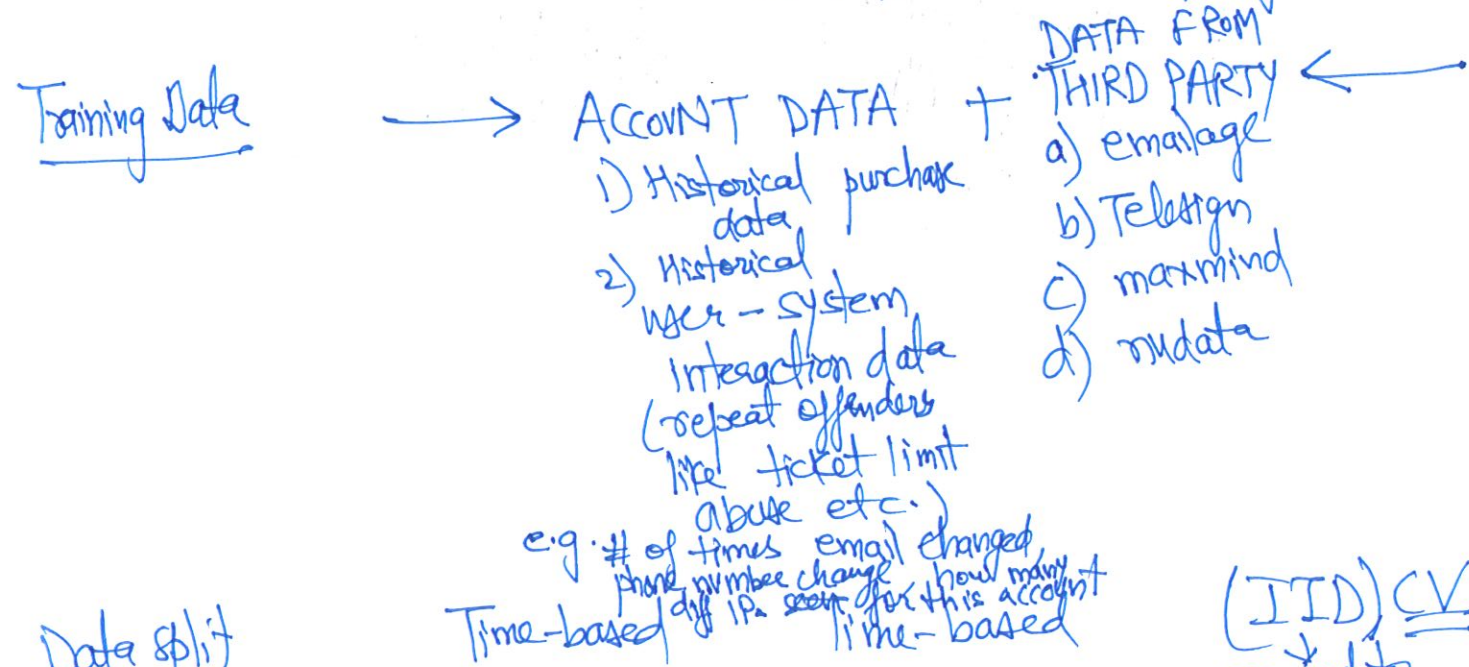
**MODELING:** — Boosting based: LightGBM (quick training on large dataset)
Fast inference
Interpretable
Re-training every week or when metrics deteriorate

DEPLOYMENT:

1) Blue green deployment (instanteous cutover but requires duplicate infra)
Shadow model testing →

ii) Canary deployment (gradual rollout to limit risk)

A/B Test :
(allows some fraud)

(only) ↙ Rule-based flagged    ↓ ↘ Model flagged (only) → cannot wait for 45-60 days to get labels.

Manual Review    |    Manual Review

| Aspect | Reserve Abuse | Fraud | IAF |
|---|---|---|---|
| Product metrics | In-cart conversion → PR | Fraud found rate → PR | Penetration metrics → PR |
| ML metrics | | | |
| Eng. | real-time/online training, real-time prediction | online prediction, offline training | batch training, batch prediction |
| Labels | ratio of reserve to purchase | Imprn → Fraud Analysts / Accredify (rule-based engine). Real labels → 45 day lag | Imprn → Listings Ind & content Market. Real labels → Presence. Real labels → only after the event has played |
| Training Data | → ACCOUNT DATA + DATA FROM THIRD PARTY ← 1) Historical purchase data 2) Historical user-system interaction data (repeat offenders like ticket limit abuse etc.) e.g. # of times email changed, phone number change, diff IPs. seen for this account | a) emailage b) Telesign c) maxmind d) nudata | |
| Data split | Time-based | Time-based | (IID) CV ↓ no need to login originally, now includes account data. |
| FEATURES (unique ones) | Velocity-based (No credit card data) | Credit card | Actor-Artist affinity (No credit card data) |

MODELS

**Vowpal Wabbit**

→ online learning

→ feature-eng. on the fly (interaction terms)

→ categorical features:
  : hashing → int
    using murmur hash
    (no hash table storage)

→ Scalable
→ Interpretable

Disad: ~~>> DYP tr~~

1) Technically linear
  → a good amt of effort in data cleaning

**LightGBM**

+ Lambda Rank Ranking for analyst review

Non-linear
(very noisy data)

→ Rule-based ensemble system had been in place:
  → moderate level of understanding of patterns

+

Tree-based models are interpretable to an extent.

**Vowpal Wabbit**

→ Roadmap for online learning

→ behavior of bad actors:
  not enough info on our problem space: hence interpretable model

1) Adv. & Disad. of Vowpal Wabbit & lightGBM
2) Some case studies in error-analysis /debugging