# Accurate, Data-Efficient Learning from Noisy, Choice-Based Labels for Inherent Risk Scoring

**W. Ronny Huang***
Ernst & Young LLP
New York, NY
ronny.huang@ey.com

**Miguel A. Perez***
Ernst & Young LLP
New York, NY
miguel.a.perez@ey.com

**Application:** Inherent risk scoring is used in anti-money laundering to determine riskiness of an entity *before* fraudulent acts occur

**Problem:** Data is scarce and the opinions of financial crime investigators are inconsistent. It is difficult to assign a risk score on an absolute scale

**Hypothesis:** We can use experts' choice-based feedback to determine the true label

## From choice to risk score

**Key point:** Obtain *absolute* information from about 5-10x the amount of *relative* information

Collected choice data encoding

$$c(y_i|l) = \begin{cases} 1, & \text{if } y_i = \max_{j \in S_k} y^j \\ -1, & \text{if } y_i = \min_{j \in S_k} y^j \\ 0, & \text{otherwise} \end{cases}$$

Average to get *mean choice*

$$\bar{c}_i(y_i) = \frac{1}{q} \sum_{l=1}^{q} c(y_i|l)$$
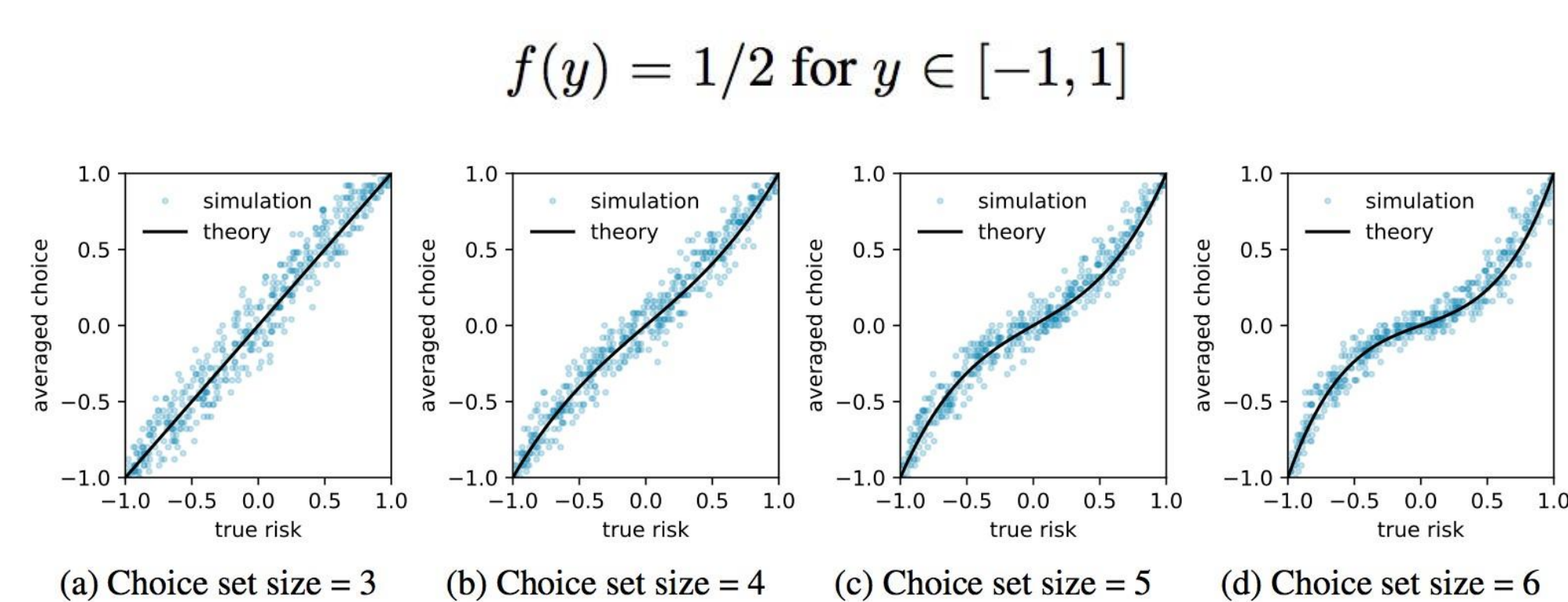
Derive the *expected choice*

$$\langle c(y_i) \rangle = \left( \int_{-\infty}^{y_i} f(y') \, \mathrm{d}y' \right)^{s-1} - \left( \int_{y_i}^{\infty} f(y') \, \mathrm{d}y' \right)^{s-1}$$

Expected choice derivation

$$\langle c(y_i) \rangle = \mathbb{E}_{y_i \sim Y} \ c(y_i)$$
$$= +1 \times P(y_i = \max_{j \in S_k} y^j)$$
$$- 1 \times P(y_i = \min_{j \in S_k} y^j)$$
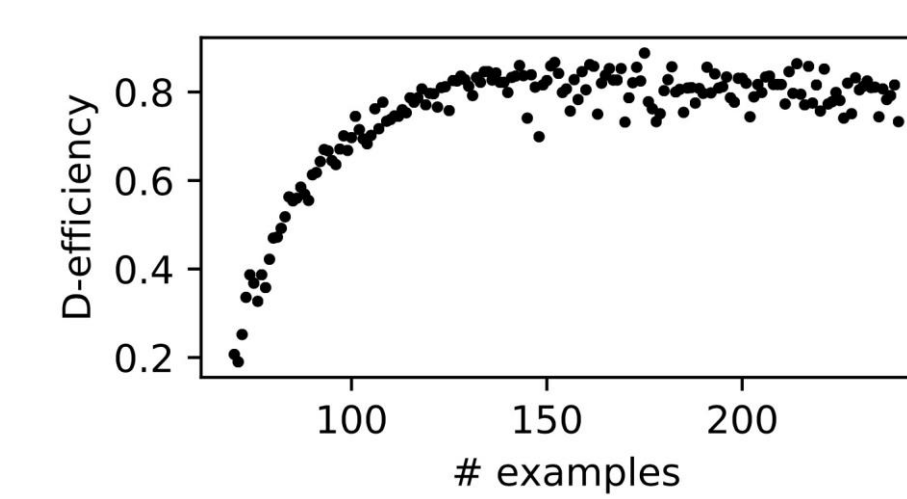$$+ 0 \times P(y_i \text{ is neither max nor min})$$

$$P(y_i = \max_{j \in S_k} y^j) = \prod_{j \in S_k, j \neq i} P(y_i > y^j)$$
$$= \prod_{j \in S_k, j \neq i} \left( \int_{-\infty}^{y_i} f(y^j) \, \mathrm{d}y^j \right)$$
$$= \left( \int_{-\infty}^{y_i} f(y') \, \mathrm{d}y' \right)^{s-1}$$

Results based on a uniform label prior

$$f(y) = 1/2 \text{ for } y \in [-1,1]$$



(a) Choice set size 3    (b) Choice set size 4    (c) Choice set size 5    (d) Choice set size 6
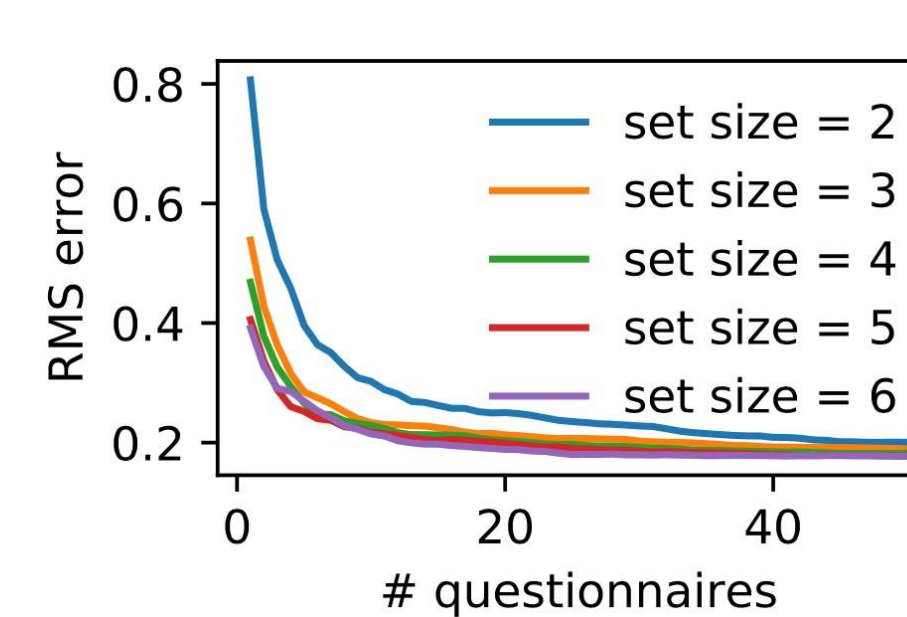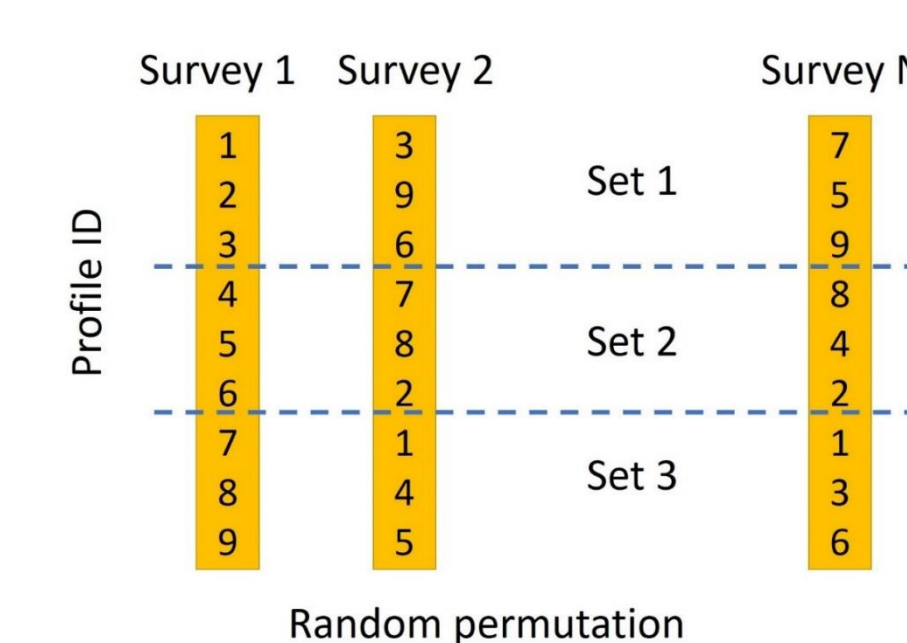
## Implementation

D-optimal synthetic dataset based on real schema



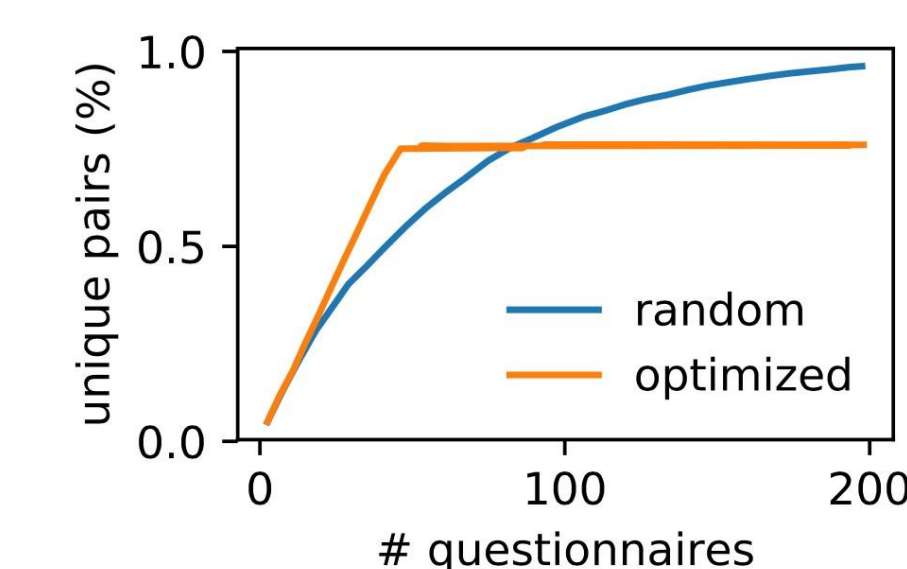Error analysis helps to determine the optimal choice set size



Choice set construction by random permutation



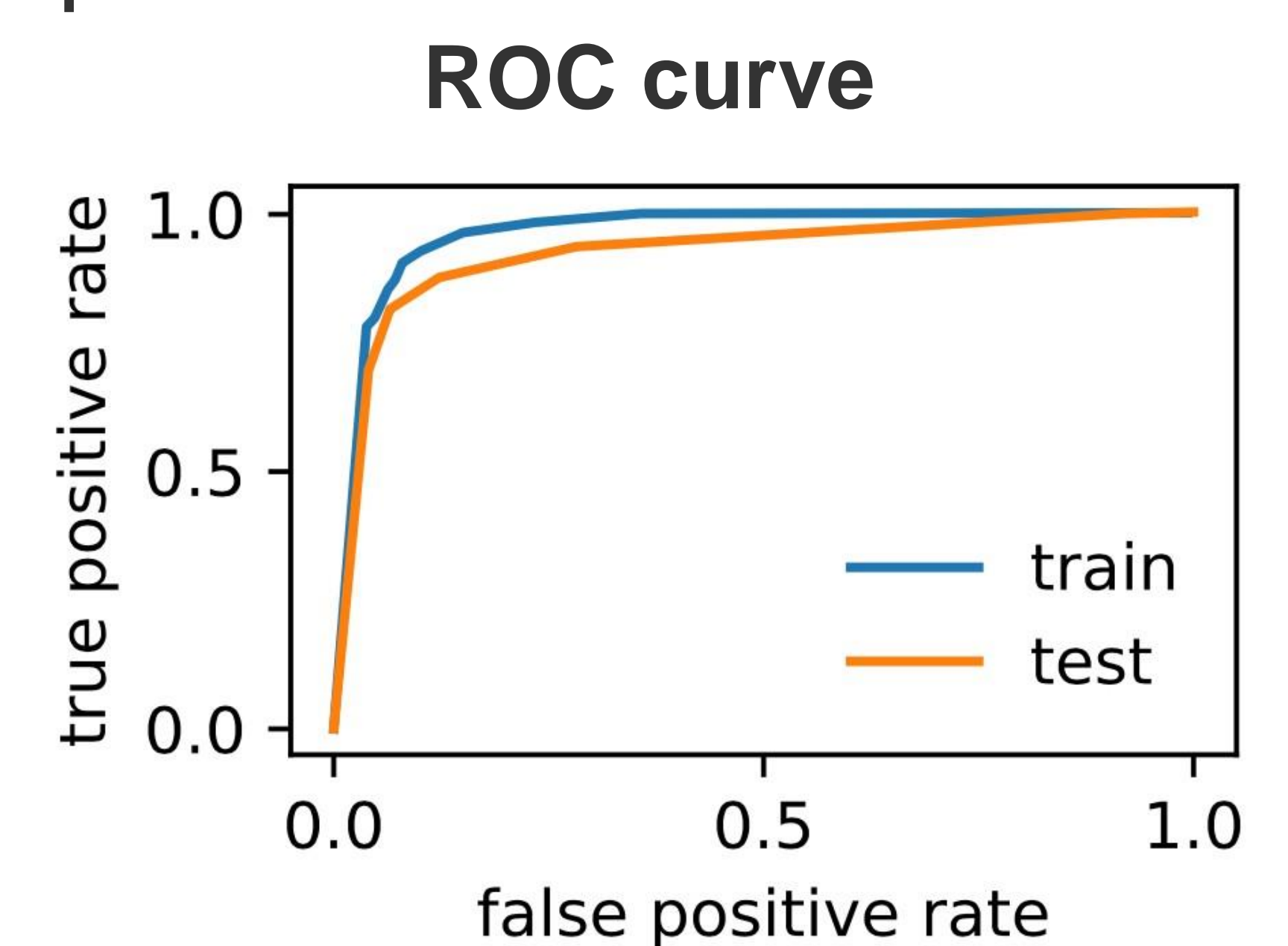Group theory enables the fastest enumeration of choice set comparisons



Choice set diversifier algorithm allows for fewer expert evaluations



## Results

**Model performance and test results at a large financial institution**

**Synthetic examples evaluated by real experts**

### ROC curve



| Metric | Train | Test |
|---|---|---|
| AUC | 97% | 93% |
| Classification Error | 8% | 11% |

### Real world examples

| Population Group | Profiles | SMA Escalations | Escalation Rate |
|---|---|---|---|
| IRM Selected Alerted Profiles | 1,500 | 28 | 1.87 |
| Remaining Scenario Alerted Profiles | 2,500 | 3 | 0.12 |

## A choice-based questionnaire for determining true labels

```
SET 1 of 25:

        age     country  pep    years
0  twenties  nigeria   no   twothree
1  thirties  nigeria  yes        one
2    fifties    italy  yes   twothree
3    sixties  nigeria  yes  fourmore
Enter indices of the most/least
risky customer (0|1|2|3): 1,0
```



(a) Example of a choice set.      (b) Example data layout.