

MOTIVATION

WANT:

- **Cheap & representative** data.
- Model charging for **difficult subgroups**.

REALITY:

- Data collection is **expensive**.
- Existing datasets are **biased**.

SOLUTION:

1. Combine existing data from **data repositories**.
2. Enforce **group count requirements**.

PROBLEM DEFINITION

GIVEN:

- **Data sources** with **sampling costs**.
- **Groups** with **minimum count requirements**.

GOAL: **Minimize** expected total **query cost**.

CONSTRAINT: Satisfy minimum count requirement.

QUERY MODEL: Uniformly random sampling.

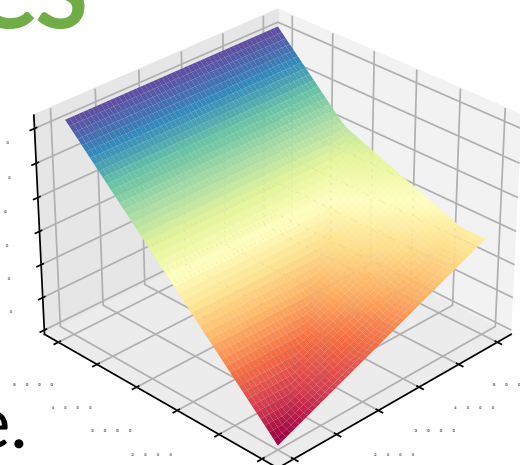
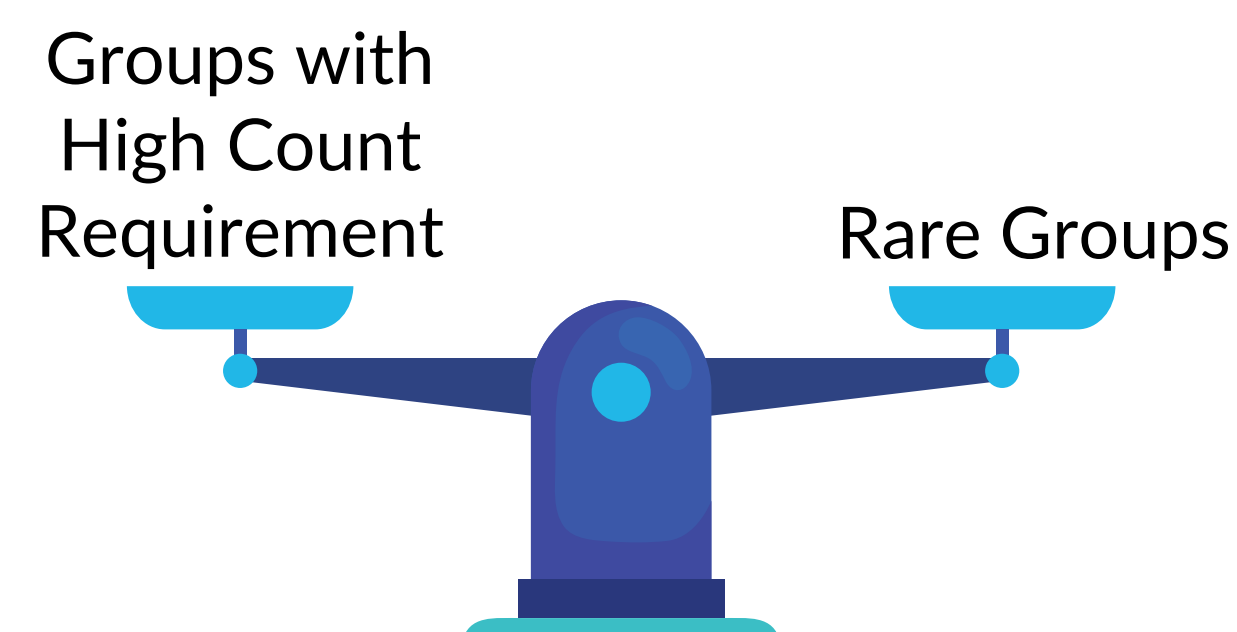
KNOWN STATISTICS

DYNAMIC PROGRAMMING

- **Optimal** but **slow** $O(Q^n)$ time.

HEURISTIC ALGORITHM **RatioColl**

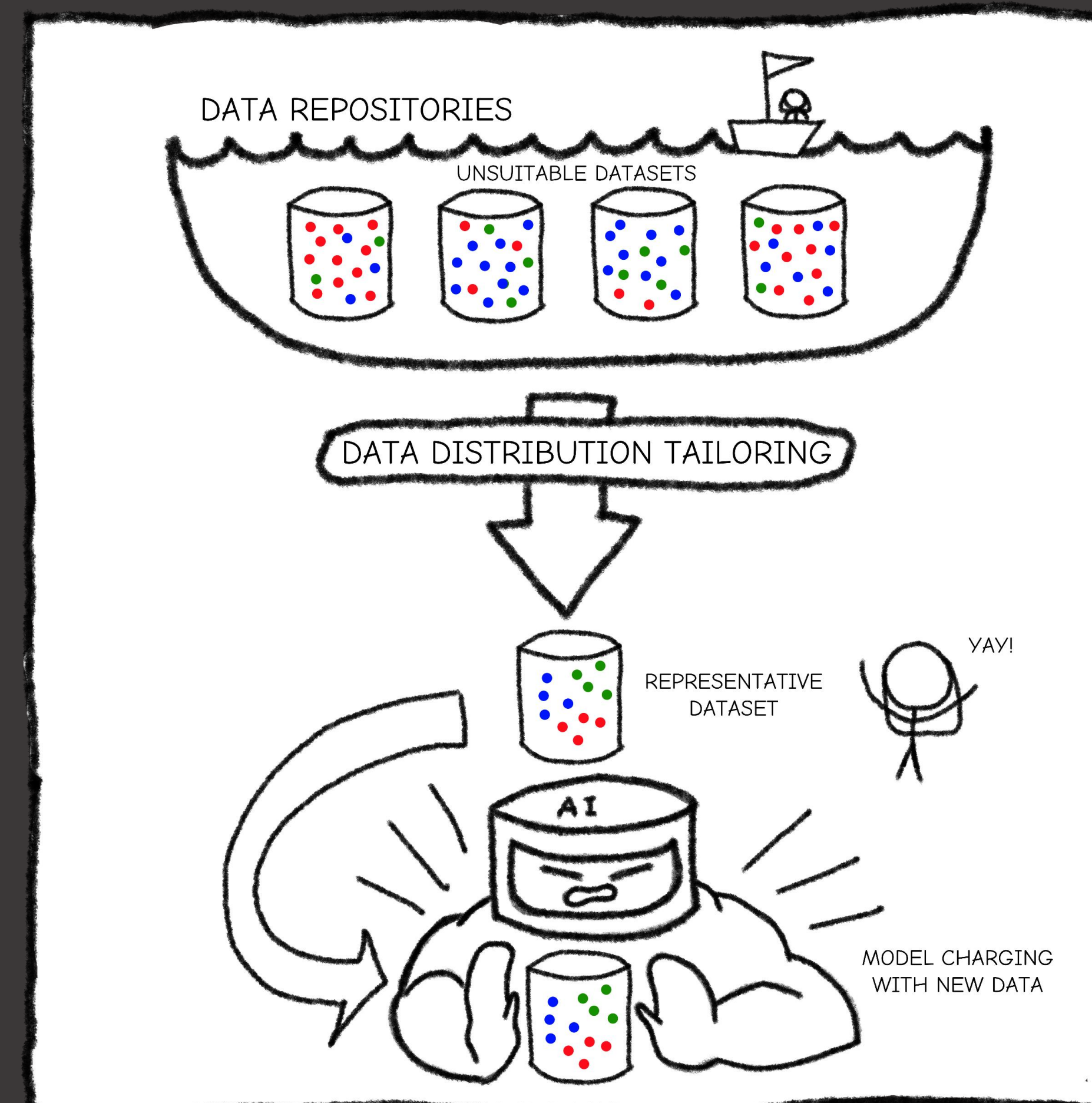
- Only \sqrt{Q} overhead cost in special case.
- **Linear** upper bound in general.



Jiwon Chang¹, Bohan Cui¹, Fatemeh Nargesian¹,
Abolfazl Asudeh², H.V. Jagadish³

¹University of Rochester ²University of Illinois Chicago ³University of Michigan

Data Distribution Tailoring Revisited: Cost-Efficient Integration of Representative Data



UNKNOWN STATISTICS

ϵ -GREEDY BANDIT

- **Same heuristic** as **RatioColl**.
- **No priors** needed.
- Sublinear **regret bound**.



ALGORITHMS

RatioColl

```

1:  $O \leftarrow \emptyset$ 
2: while query is not satisfied:
3:    $G^* = \operatorname{argmax}_{G_j, Q_j > 0} \left( Q_j \cdot \min_{i \in [n]} \left( \frac{c_i}{p_{i,j}} \right) \right)$  // Choose priority group
4:    $D^* = \operatorname{argmin}_{D_i} \left( \frac{c_i}{p_{i,*}} \right)$  // Maximize priority group per cost
5:    $s \leftarrow \text{Query}(D^*)$ 
6:    $O \leftarrow O \cup \{s\}$ 
7: return  $O$ 

```

EpsilonGreedy

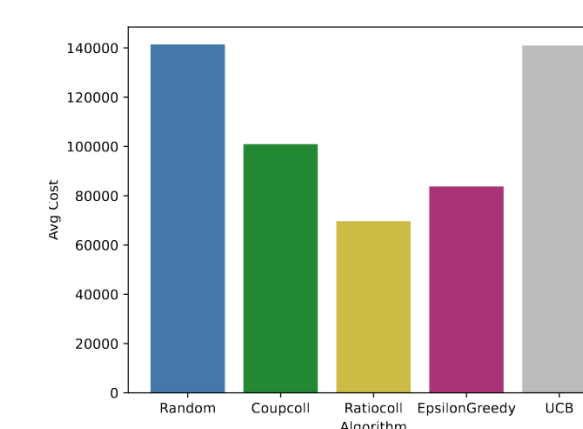
```

1:  $O \leftarrow \emptyset$ 
2: while query is not satisfied:
3:   if  $t \leq n$  then  $D^* \leftarrow D_t$  // Initialization
4:   else:
5:     with probability  $p = \sqrt[3]{\ln t / t}$ : // Exploration round
6:        $D^* \leftarrow$  random data source
7:     else: // Greedy exploitation round
8:       let  $\bar{p}_{i,j}$   $\leftarrow$  estimate of  $p_{i,j}$  based on sample mean
8:        $R(G_j) \leftarrow \left( Q_j \cdot \min_{i \in [n]} \left( \frac{c_i}{\bar{p}_{i,j}} \right) \right)$  for each  $G_j$  // Reward of group
9:        $D^* \leftarrow \operatorname{argmax}_{D_i} \left( \frac{1}{c_i} \sum_{j \in [m]} \bar{p}_{i,j} \cdot R(G_j) \right)$  // Reward of data source
10:   $s \leftarrow \text{Query}(D^*)$ 
11:   $O \leftarrow O \cup \{s\}$ 
12:  update trackers for probability estimation
13: return  $O$ 

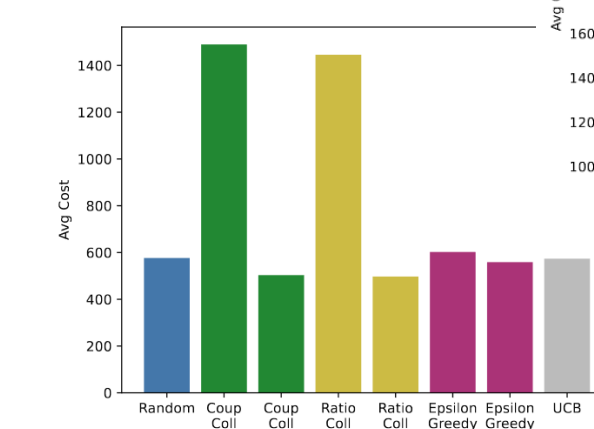
```

SELECTED RESULTS

- RatioColl consistently out-performs SOTA.
- EpsilonGreedy competitive with SOTA despite needing no priors.



Flights Dataset



COMPAS

