

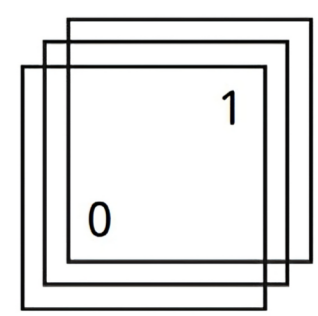


Diverse Data Selection under Fairness Constraints

Zafeiria (Iro) Moumoulidou
zmoumoulidou@cs.umass.edu

Andrew McGregor
mcgregor@cs.umass.edu

Alexandra Meliou
ameli@cs.umass.edu



DREAM LAB

DATA SYSTEMS RESEARCH FOR EXPLORATION,
ANALYTICS, AND MODELING

The Fair-Swap algorithm: # colors = 2

Universe of points placed on a line:

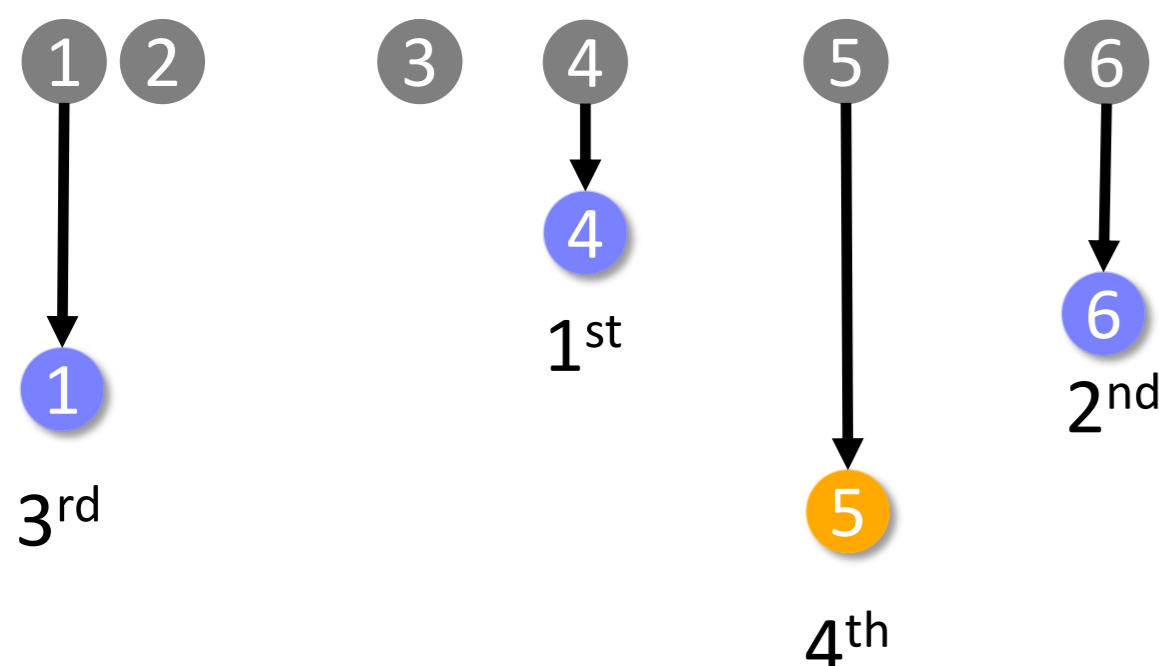


Select a set of 4 points with:

blue = 2 # orange = 2

A. Color-blind Phase

1. Use farthest-first traversal heuristic to pick the points.



There is a +1 blue point and a -1 orange point is missing.
Let's **swap**!

B. Swap Phase

1. Let's **add** an **orange** point!



We add point 2 because it's the **farthest away** from the other orange point 5!

2. Let's **remove** a **blue** point!



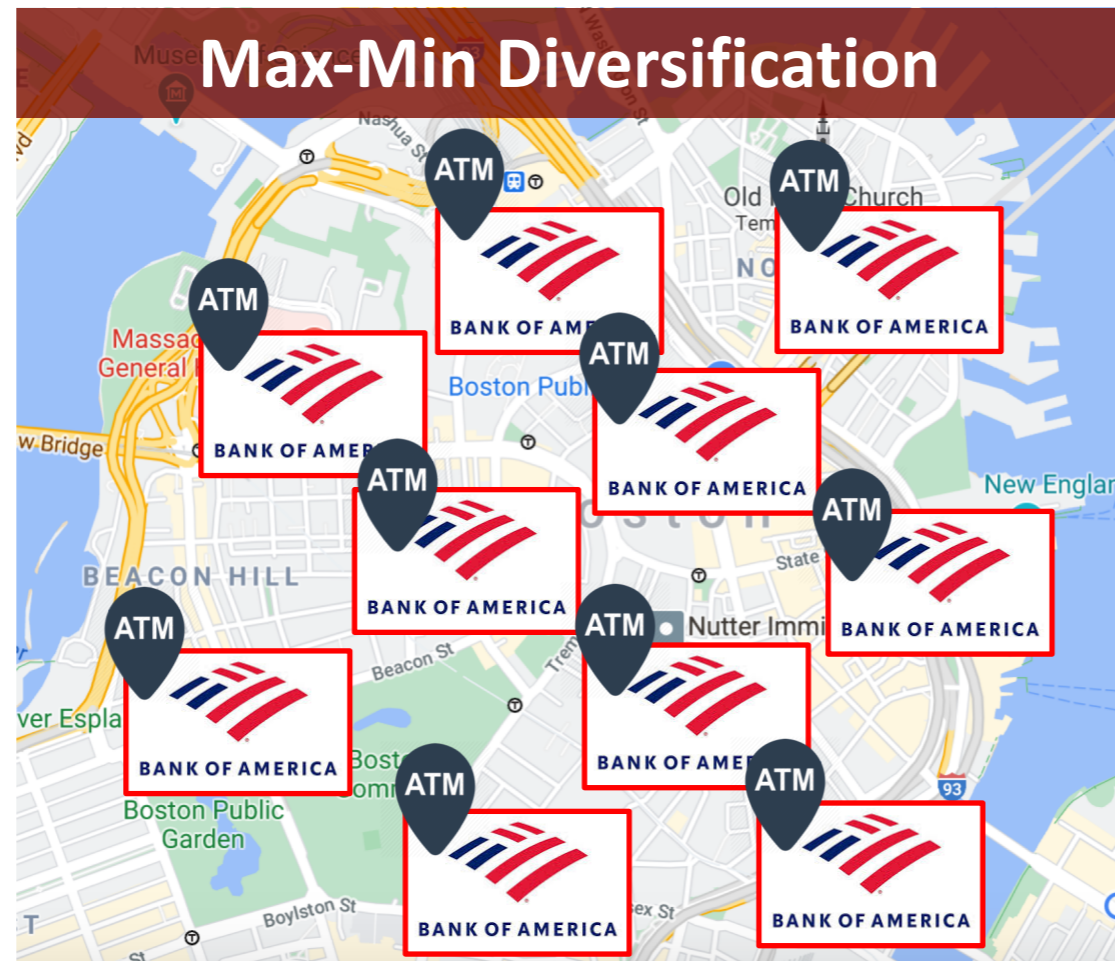
We remove point 1 because it's the **closest** to point 2 we just added!

We are done!

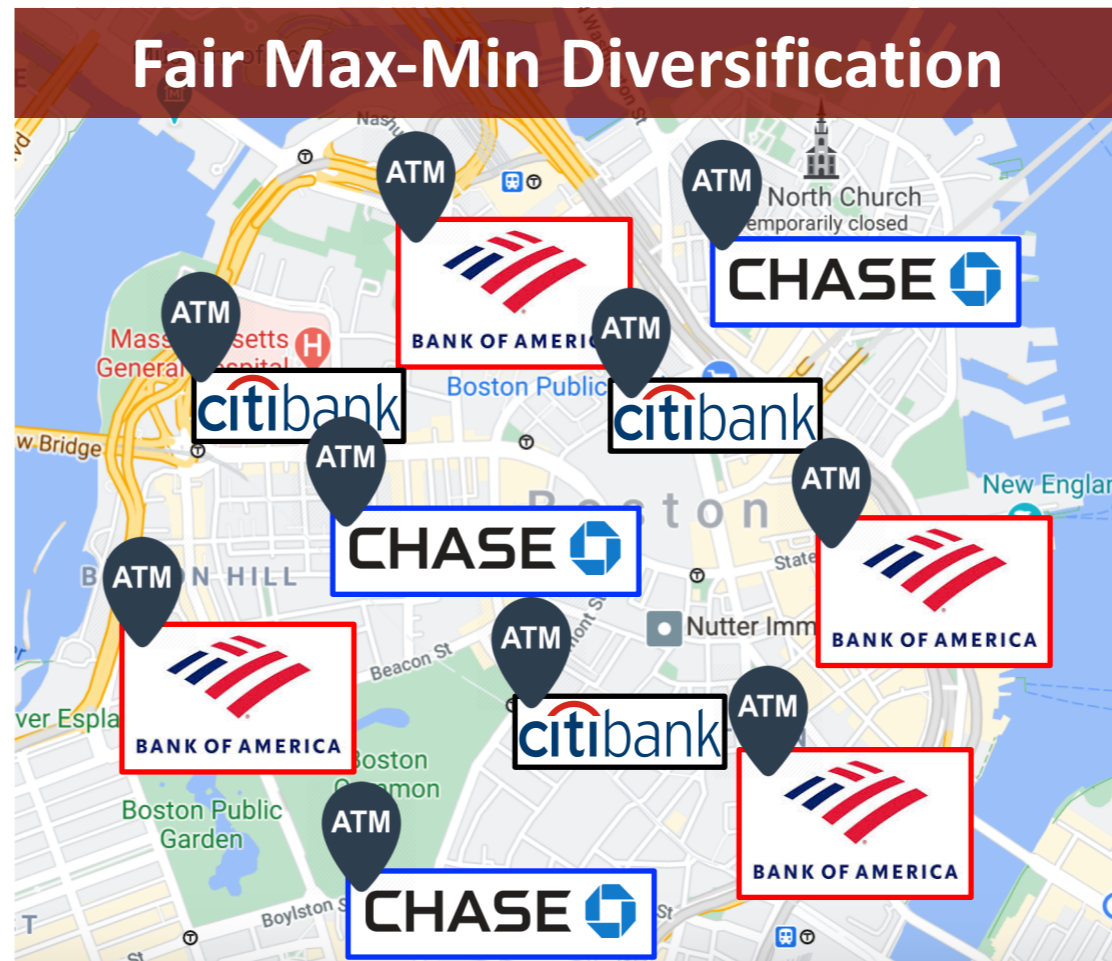


Fair and Diverse Data Selection

Max-Min Diversification



Fair Max-Min Diversification



The Fair Max-Min diversification problem



FAIR MAX-MIN : maximize $\min_{\substack{u,v \in S \\ u \neq v}} d(u,v)$
NP-hard
subject to $|S \cap \mathcal{U}_i| = k_i, \forall i \in [m]$
 $m = \# \text{ colors}$

Can we do better?

Moumoulidou et al.
[ICDT 2021]

4— approx. ($m = 2$)
(3m-1)— approx. ($m \geq 3$)

Addanki et al.
[ICDT 2022]



$m = \# \text{ colors}$ $\epsilon \in [0,1]$: approx. error parameter

General Metric Spaces

2— approx. with expected fairness
6— approx. with $(1-\epsilon)$ fairness

Euclidean Metric Spaces

Exact solution. (1-dimensional spaces)
(1+ ϵ)— approx. with $(1-\epsilon)$ fairness
+ Streaming & Distributed Implementations

The Fair-Flow algorithm: # colors ≥ 3

Universe of points placed on a line:



Select a set of 4 points with:

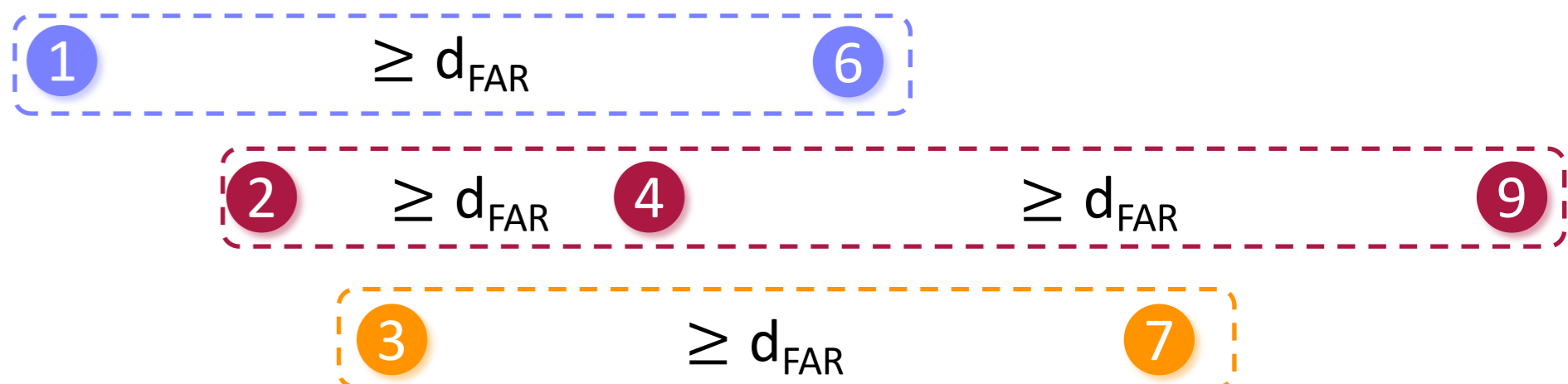
blue = 1 # orange = 1 # purple = 2

Oracle d_{FAR}
 d_{CLOSE}

A. Step 1

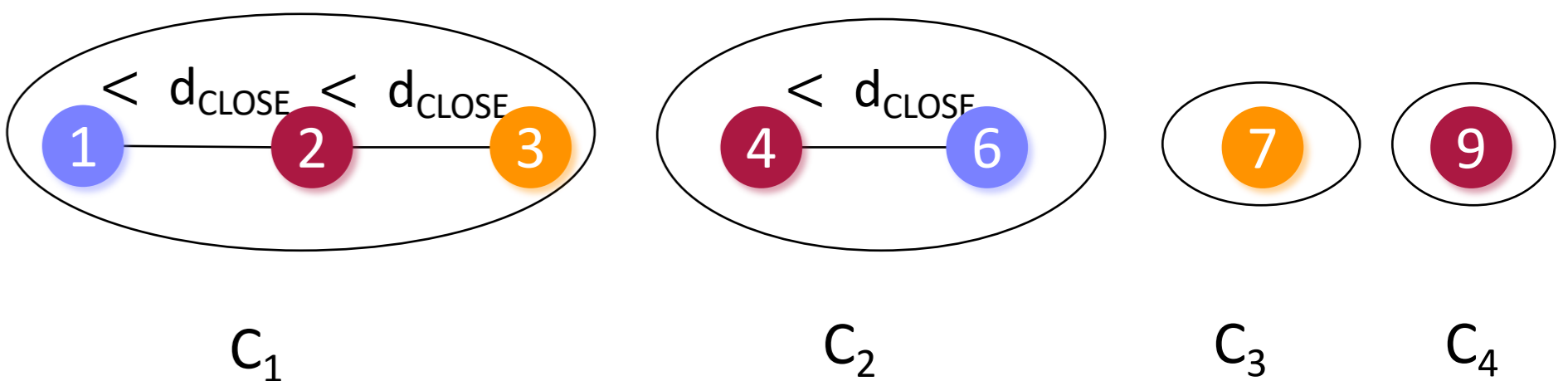
For every color:

Find a maximal set of points that are $\geq d_{\text{FAR}}$ far apart.



B. Step 2

Connect with an edge any two points $< d_{\text{CLOSE}}$ far apart.



C. Step 3

Solve a Max-Flow problem to find a fair and diverse set.

