

CSCSE 689: Special Topics in Modern Algorithms for Data Science

Lecture 34

Samson Zhou

Presentation Schedule

- **November 27:** Chunkai, Jung, Galaxy AI
- **November 29:** STMI, Anmol, Jason
- **December 1:** Bokun, Ayesha, Dawei, Lipai

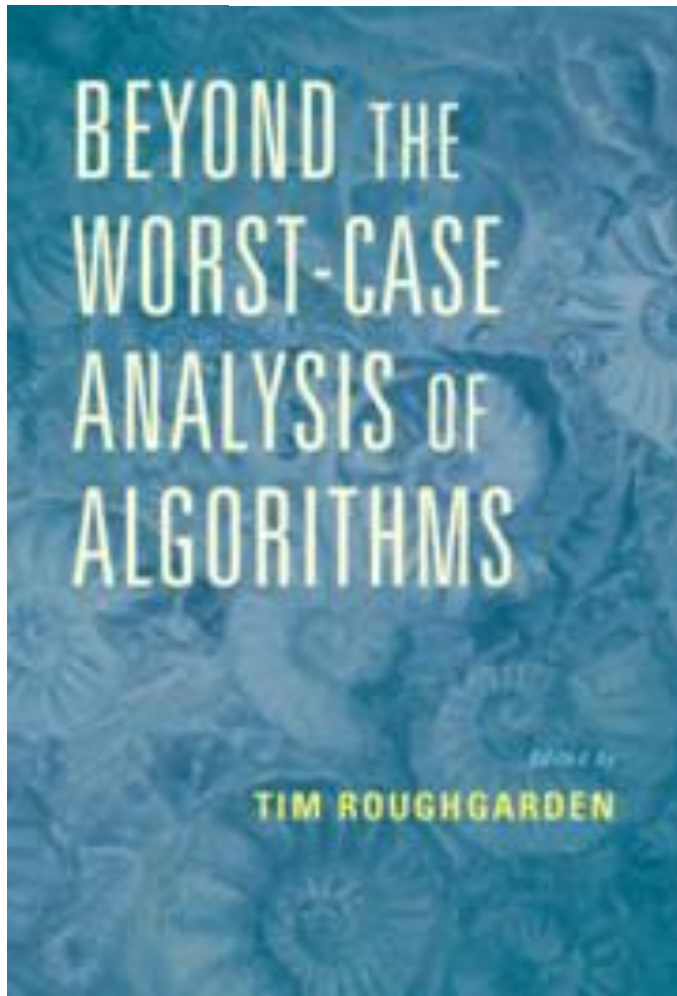


Stellar™

6.890 Learning-Augmented Algorithms

LOGIN

Course : » Course 6 : » Spring 2019 : » 6.890 : » Materials



BEYOND THE WORST-CASE ANALYSIS OF ALGORITHMS

EDITED BY
TIM ROUGHGARDEN

29	Data-Driven Algorithm Design	626
	<i>Maria-Florina Balcan</i>	
29.1	Motivation and Context	626
29.2	Data-Driven Algorithm Design via Statistical Learning	628
29.3	Data-Driven Algorithm Design via Online Learning	639
29.4	Summary and Discussion	644
30	Algorithms with Predictions	646
	<i>Michael Mitzenmacher and Sergei Vassilvitskii</i>	
30.1	Introduction	646
30.2	Counting Sketches	649
30.3	Learned Bloom Filters	650
30.4	Caching with Predictions	652
30.5	Scheduling with Predictions	655
30.6	Notes	660

Learning-Augmented Algorithms

- For a certain task and input, algorithm is given advice
- Advice could be “good”, advice could be “bad”
- **Goal:** “Good” performance if the advice is good, “normal” performance if the advice is bad



Mahabalipuram

73°

Partly Cloudy
H:86° L:72°

Sunny conditions expected around 9AM.

Now	8AM	9AM	10AM	11AM	12P
73°	76°	79°	82°	84°	85°

10-DAY FORECAST

Today		72°	86°
Fri		72°	86°
Sat		71°	85°
Sun		71°	86°

Mahabalipuram

77°

Partly Cloudy
H:84° L:76°

Sunny conditions expected around 9AM.

Now	8AM	9AM	10AM	11AM	12P
77°	78°	80°	82°	83°	83°

10-DAY FORECAST

Today		76°	84°
Fri		75°	84°
Sat		74°	85°
Sun		74°	85°
Mon		71°	85°

Mahabalipuram

74°

Partly Cloudy
H:85° L:73°

Sunny conditions expected around 9AM.

Now	8AM	9AM	10AM	11AM	12PM
74°	76°	79°	81°	83°	84°

10-DAY FORECAST

Today		73°	85°
Fri		73°	85°
Sat		72°	84°
Sun		72°	84°

Learning-Augmented Algorithms

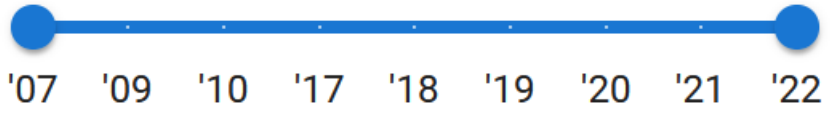
- **Better data structures:** Bloom filters with lower false positive rates [Mitzenmacher18]
- **Better space-accuracy tradeoff for streaming algorithms:** Frequency estimation, e.g., CountMin, CountSketch [HsuIndykKatabiVakilian19], moment estimation, distinct elements [JiangLinRuanWoodruff20], triangle counting [ChenEdenIndykLinNarayananRubinfeldSilwalWagnerWoodruffZhang22]
- **Better size-accuracy tradeoff for sketching:** Low-rank approximation [IndykVakilianYuan19]

Learning-Augmented Algorithms

- **Warm-start to search algorithms:** Binary search [LinLuoWoodruff22], Max-flow [ChenSilwalVakilianZhang22], [DaviesMoseleyVassilvitskiiWang23], matchings [DinitzImLavastidaMoseleyVassilvitskii21]
- **Better accuracy-sample complexity tradeoff:** Support size estimation [EdenIndykNarayananRubinfeldSilwalWagner21]
- **Better online algorithms:** Set cover [BamasMaggioriSvensson20], [GrigorescuLinSilwalSongZhou23], Scheduling [LattanziLavastidaMoseleyVassilvitskii20], [ScullyGrosfMitzenmacher22]
- **Better privacy-utility tradeoffs for DP:** Quantile estimation [KhodakAminDickVassilvitskii23]
- **Beating NP-hardness?**

Algorithms with Predictions

PAPER LIST FURTHER MATERIAL ABOUT



Newest first ▾

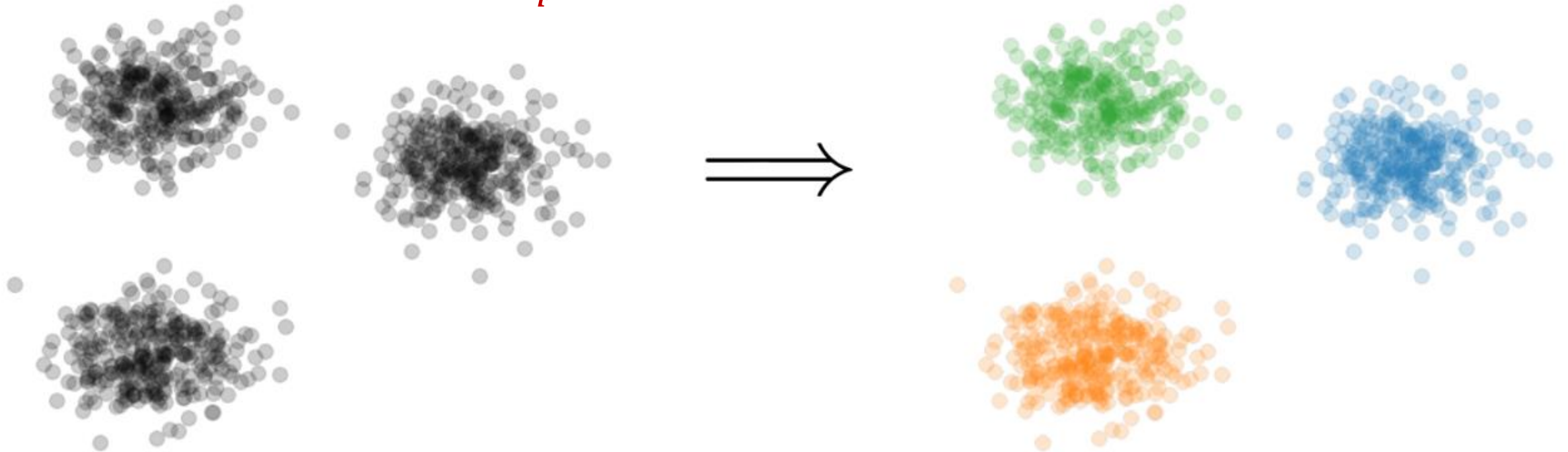
122 papers

- Graph Searching with Predictions** Banerjee, Cohen-Addad, Gupta, Li arXiv '22 exploration online search
- Scheduling with Predictions** Cho, Henderson, Shmoys arXiv '22 online scheduling
- On the Power of Learning-Augmented BSTs** Chen, Chen arXiv '22 data structure search
- Algorithms with Prediction Portfolios** Dinitz, Im, Lavastida, Moseley, Vassilvitskii arXiv '22 load balancing matching multiple predictions online scheduling
- Private Algorithms with Private Predictions** Amin, Dick, Khodak, Vassilvitskii arXiv '22 differential privacy
- Paging with Succinct Predictions** Antoniadis, Boyar, Eliáš, Favrholdt, Hoeksma, Larsen, Polak, Simon arXiv '22 caching/paging online
- Proportionally Fair Online Allocation of Public Goods with Predictions** Banerjee, Gkatzelis, Hossain, Jin, Micha, Shah arXiv '22 allocation online
- Canadian Traveller Problem with Predictions** Bampis, Escoffier, Xeferis arXiv '22 WAOA '22 online routing
- Learning-Augmented Algorithms for Online Linear and Semidefinite Programming** Grigorescu, Lin, Silwal, Song, Zhou arXiv '22 covering problems online SDP

Learning-Augmented Clustering

- **Goal:** Given dataset P in d dimensions, output a set C of k centers to minimize

$$\sum_{p \in P} \min_{c \in C} \|p - c\|_2^2$$



Learning-Augmented Clustering

- **Goal:** Given dataset P in d dimensions, output a set C of k centers to minimize

$$\sum_{p \in P} \min_{c \in C} \|p - c\|_2^2$$

- **NP-hard** to even approximate within a factor of **1.07** [Cohen-AddadC.S.20, LeeSchmidtWright17]
- **Beyond worst-case:** Clustering on inputs from some “nice” distribution, similar inputs or inputs with auxiliary information
- **Hope:** ML can guide the clustering, so we can overcome worst-case with advice!

Predictor

- Suppose Π outputs noisy labels according to a $(1 + \alpha)$ approximate clustering C and error rate $\lambda \leq \alpha$



What is the label of x_1 ?

What is the label of x_2 ?

x_1 belongs to cluster 3

x_2 belongs to cluster 7



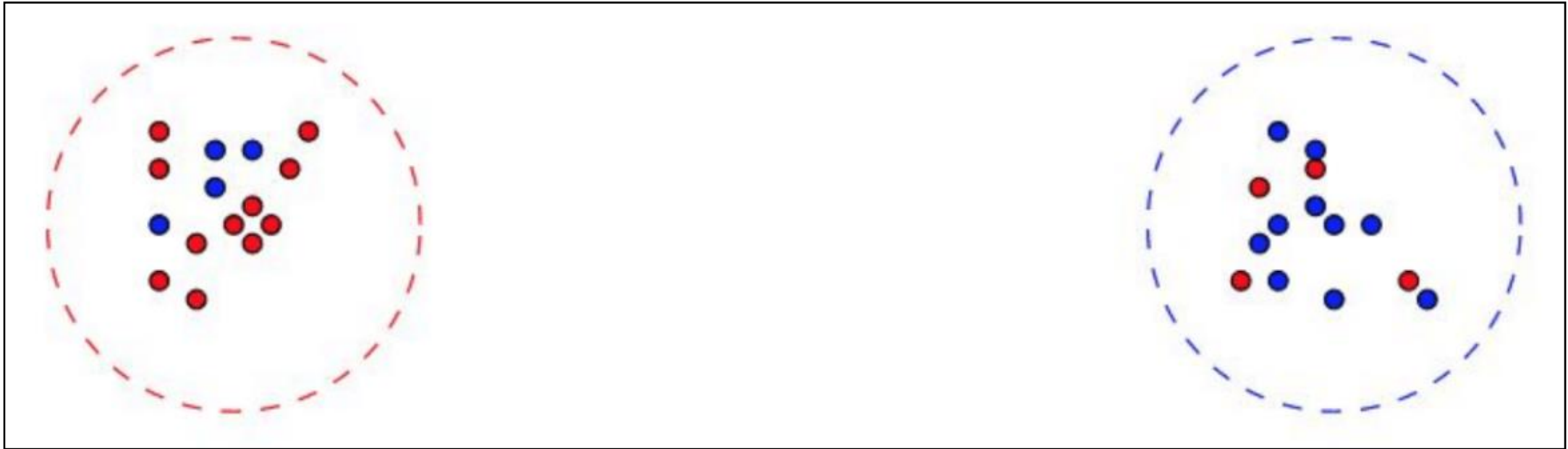
Theoretical Guarantee

- Suppose Π outputs noisy labels according to a $(1 + \alpha)$ approximate clustering C and error rate $\lambda \leq \alpha$
- **Main result [EFSWZ22]**: Algorithm that outputs a $(1 + O(\alpha))$ approximate k -means clustering in nearly linear time

- “Predictions can overcome complexity hardness barriers!”

Naïve Approach Does Not Work

- Not enough to blindly follow predictions!



- Optimal cost ≈ 0
- Predictor with arbitrary small error has large cost!

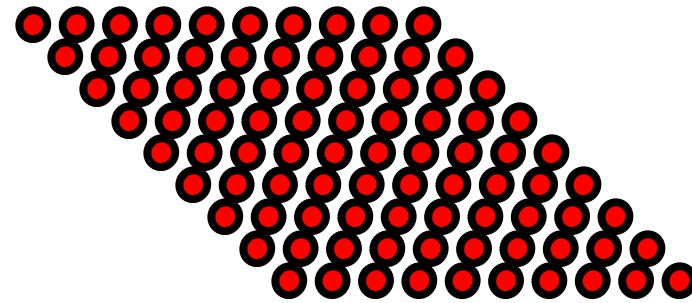
Naïve Approach Does Not Work

- Can a predictor even help?

Cluster 1

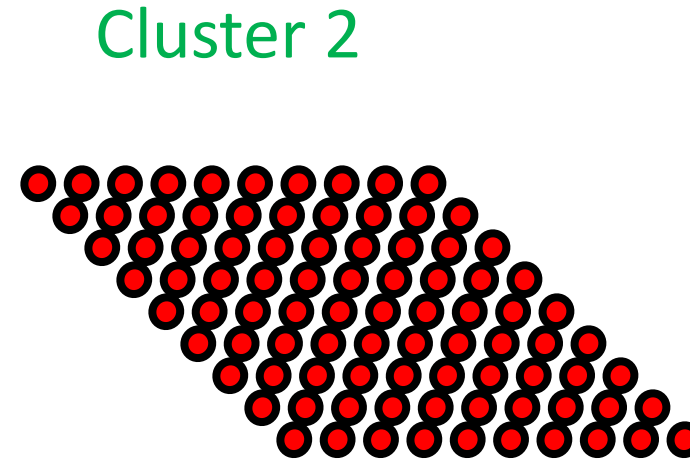
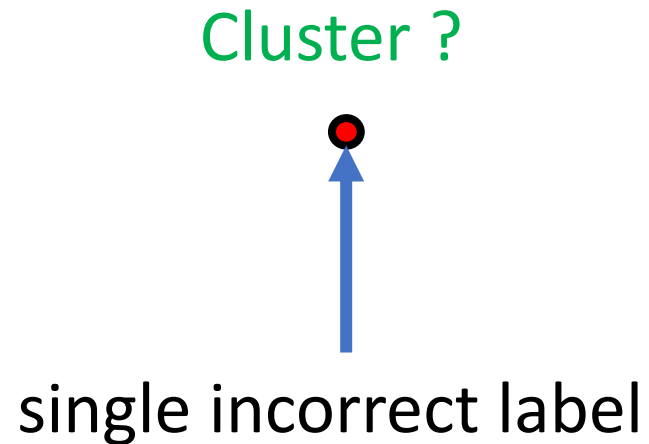


Cluster 2



Naïve Approach Does Not Work

- Can a predictor even help?



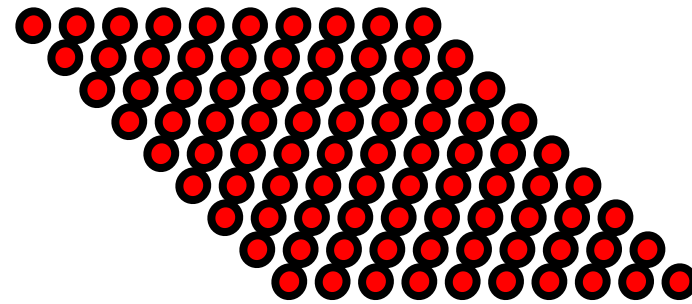
Naïve Approach Does Not Work

- Can a predictor even help?

Cluster ?



Cluster 2




- **MUST** have assumptions about the accuracy on each cluster


Precision and Recall

- [EFSWZ22]: Assume cluster sizes are “balanced”
- [NCN23]: Let P_i be the optimal cluster with label i and Q_i be the points that are labeled i . Then $|Q_i \setminus P_i| + |P_i \setminus Q_i| \leq \alpha \cdot |P_i|$.

Precision

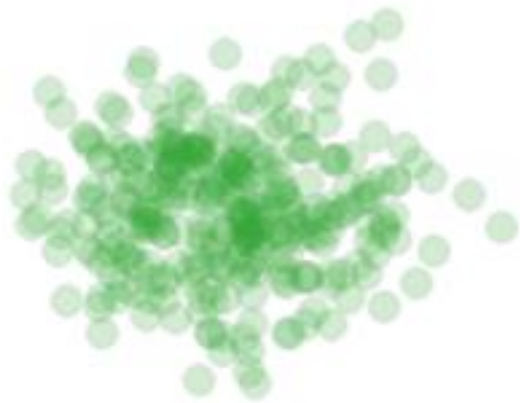


Recall



Algorithmic Intuition

- **Our approach:** Closed-form solution for best center of a *fixed* set of points



$$\operatorname{argmin}_c [\operatorname{cost}(c, P)] = \frac{1}{|P|} \sum_{p \in P} p$$

$$\operatorname{argmin}_c \sum_{p \in P} \|p - c\|_2^2 = \frac{1}{|P|} \sum_{p \in P} p$$

Algorithmic Intuition

- Consider each *dimension* separately

Algorithm 1 Learning-augmented k -means clustering

Input: A point set X with labels given by a predictor Π with label error rate λ

Output: $(1+O(\alpha))$ -approximate k -means clustering of X

- 1: **for** $i = 1$ to $i = k$ **do**
 - 2: Let Y_i be the set of points with label i .
 - 3: Run CRDEST for each of the d coordinates of Y_i .
 - 4: Let C'_i be the coordinate-wise outputs of CRDEST.
 - 5: **end for**
 - 6: **Return** C'_1, \dots, C'_k .
-

Algorithmic Intuition

- Consider each *label* separately

Algorithm 1 Learning-augmented k -means clustering

Input: A point set X with labels given by a predictor Π with label error rate λ

Output: $(1+O(\alpha))$ -approximate k -means clustering of X

1: **for** $i = 1$ to $i = k$ **do**

2: Let Y_i be the set of points with label i .

3: Run CRDEST for each of the d coordinates of Y_i .

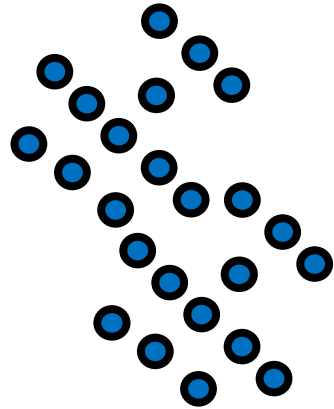
4: Let C'_i be the coordinate-wise outputs of CRDEST.

5: **end for**

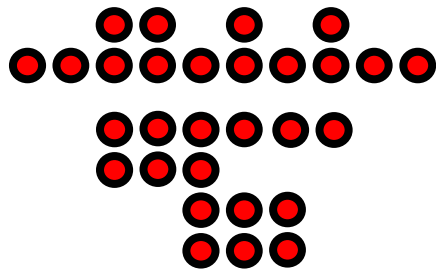
6: **Return** C'_1, \dots, C'_k .

Algorithmic Intuition

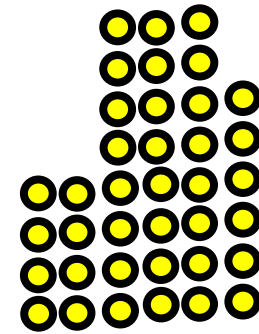
- Example:



Cluster 2



Cluster 1



Cluster 3

Algorithmic Intuition

- **Example:** Consider the points with label **1**



Algorithmic Intuition

- **Example:** Consider the points with label **1**
- Consider each dimension separately



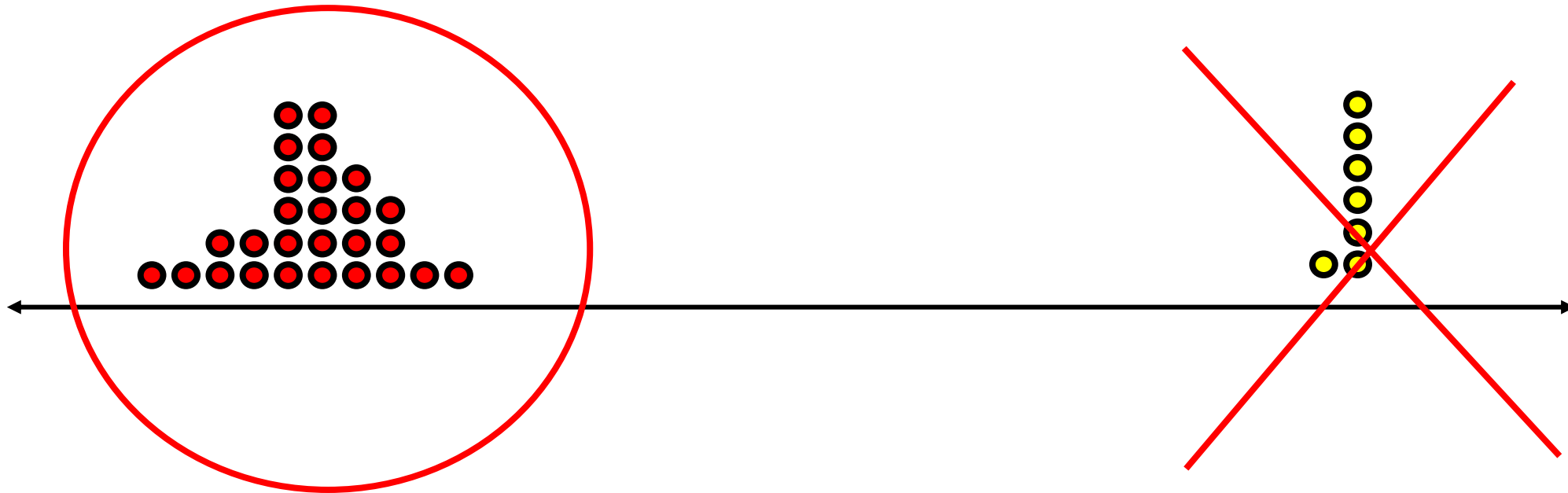
Algorithmic Intuition

- **Example:** Consider the histogram of points with label **1**



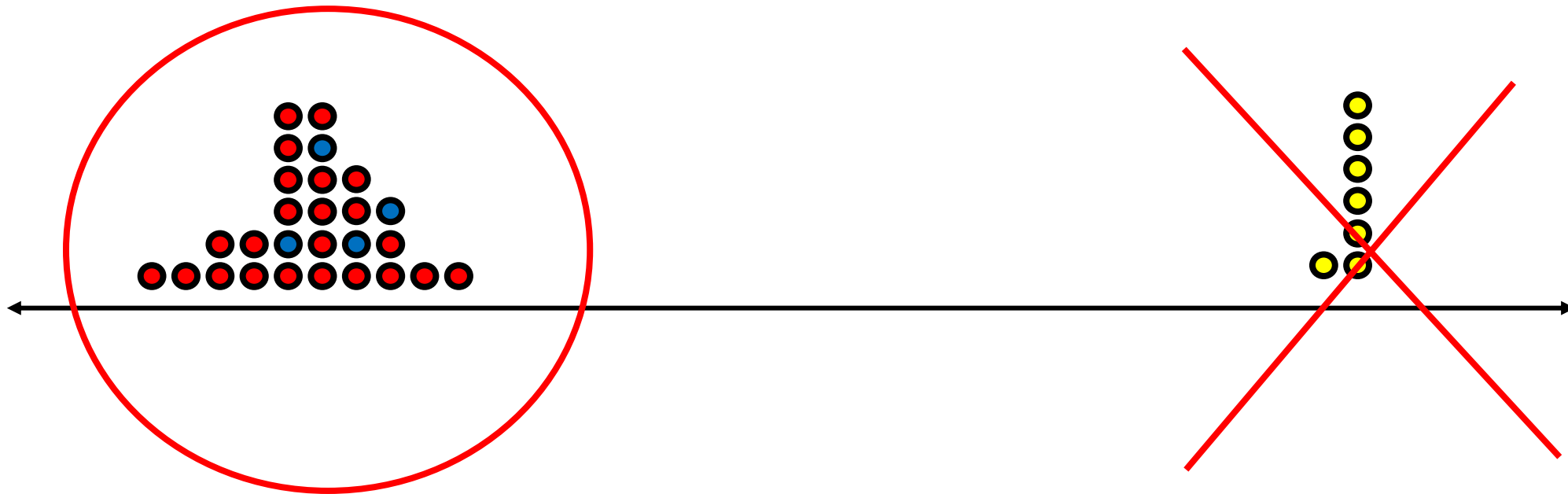
Algorithmic Intuition

- **Example:** Consider the histogram of points with label **1**
- Is it true that “pruning” away the outliers removes all incorrect points?



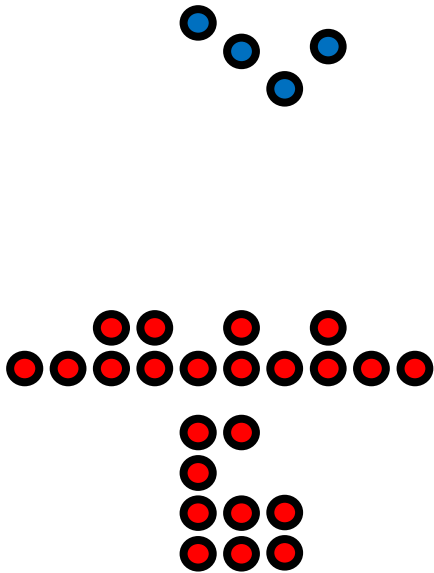
Algorithmic Intuition

- Is it true that “pruning” away the outliers removes all incorrect points? **NO!**



Algorithmic Intuition

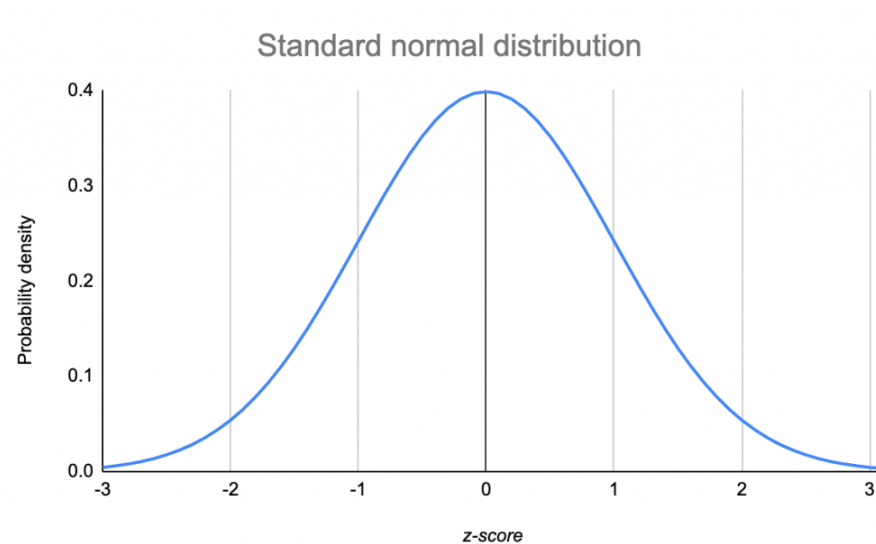
- **Example:** Consider the points with label **1**



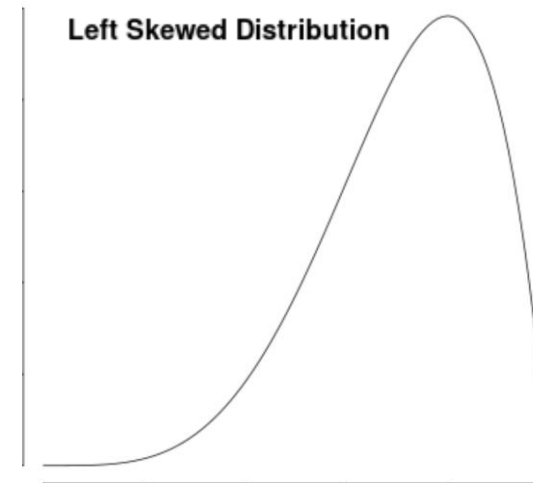
Algorithmic Intuition

- Consider each label and each dimension separately
- **Our approach:** Use ideas from robust mean estimation

$$(1 - \alpha)P$$



$$\alpha Q$$



Algorithmic Intuition

- Case 1: Q is “far” from P



Algorithmic Intuition

- Case 1: Q is “far” from P
- Can detect handle this case by “pruning” the distribution



Algorithmic Intuition

- Case 2: Q is “close” to P



Algorithmic Intuition

- Case 2: Q is “close” to P
- Q cannot heavily affect the empirical mean P



Algorithmic Intuition

- **Algorithm:** Find the mean of the shortest interval that contains $(1 - O(\alpha))$ fraction of the points



Algorithm

- **Algorithm:** Find the mean of the shortest interval that contains $(1 - O(\alpha))$ fraction of the points

Algorithm 2 Coordinate-wise estimation CRDEST

Input: Points $x_1, \dots, x_{2m} \in \mathbb{R}$, corruption level $\lambda \leq \alpha$

- 1: Randomly partition the points into two groups X_1, X_2 of size m .
 - 2: Let $I = [a, b]$ be the shortest interval containing $m(1 - 5\alpha)$ points of X_1 .
 - 3: $Z \leftarrow X_2 \cap I$
 - 4: $z \leftarrow \frac{1}{|Z|} \sum_{x \in Z} x$
 - 5: **Return** z
-

Analysis Overview

- Robust mean estimation gives additive α error to the *location* of the mean
- How does this affect the k -means clustering cost?

Analysis Overview

- **Analysis:** Robust mean gives $(1 + \alpha)$ -approximation to the **1**-means clustering cost
- **Recall:** Consider each label and each dimension separately



Analysis Overview

- **Analysis:** Robust mean gives $(1 + \alpha)$ -approximation to the k -means clustering cost
- **Lemma:** Let P, Q be sets of real numbers with $|P| \geq (1 - \alpha)n$ and $|Q| \leq \alpha n$. Let $X = P \cup Q$, let C_X and C_P be the means of X and P . Then

$$\text{Cost}(X, C_P) \leq (1 + \alpha)\text{Cost}(X, C_X)$$

- [InabaKatoHImai94]:

$$\text{Cost}(X, C_P) \leq \text{Cost}(X, C_X) + |X| \cdot |C_P - C_X|^2$$

Algorithm 1 Learning-augmented k -means clustering

Input: A point set X with labels given by a predictor Π with label error rate λ

Output: $(1+O(\alpha))$ -approximate k -means clustering of X

- 1: **for** $i = 1$ to $i = k$ **do**
 - 2: Let Y_i be the set of points with label i .
 - 3: Run CRDEST for each of the d coordinates of Y_i .
 - 4: Let C'_i be the coordinate-wise outputs of CRDEST.
 - 5: **end for**
 - 6: **Return** C'_1, \dots, C'_k .
-

Algorithm 2 Coordinate-wise estimation CRDEST

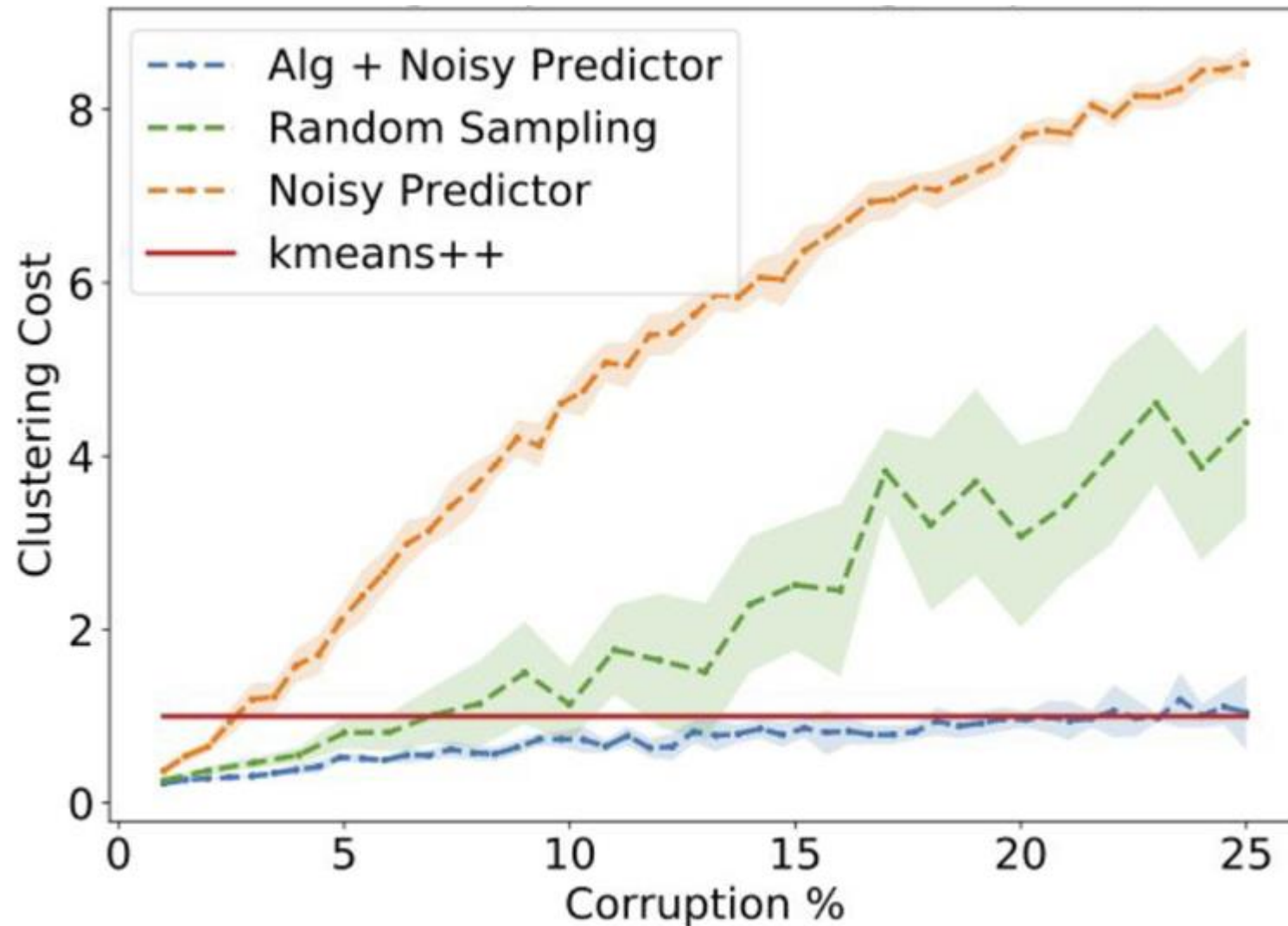
Input: Points $x_1, \dots, x_{2m} \in \mathbb{R}$, corruption level $\lambda \leq \alpha$

- 1: Randomly partition the points into two groups X_1, X_2 of size m .
 - 2: Let $I = [a, b]$ be the shortest interval containing $m(1 - 5\alpha)$ points of X_1 .
 - 3: $Z \leftarrow X_2 \cap I$
 - 4: $z \leftarrow \frac{1}{|Z|} \sum_{x \in Z} x$
 - 5: **Return** z
-

Experimental Results

- **Case Study:** Spectral clustering on graphs varying over time
- **Dataset:** Internet router graph varying over the course of a year
- **Methodology:** Compare to standard benchmarks while using various natural predictors, i.e., noisily perturb true labels and compare to baselines as function of error

Dataset: Internet router graph varying over the course of a year, $k = 10$



Conclusion: Our algorithm (using predictor) outperforms benchmarks such as k -means ++ for low error while staying competitive with high corruptions

Summary

- **NP-hard** to even approximate within a factor of **1.07** [Cohen-AddadC.S.20, LeeSchmidtWright17]
- **Main result** [EFSWZ22]: Algorithm that outputs a $(1 + O(\alpha))$ approximate k -means clustering in nearly linear time
- Handles clustering with *outliers*
- Not enough to blindly follow predictions!
- **Our approach**: Use ideas from robust mean estimation

...and Beyond!

- **Related work:**
- Semi-supervised active clustering (SSAC) framework: Same cluster queries, [AKB16], [KG17], [MS17], [GHS18], [ABJK18], ..., correlation clustering
- **Future directions:**
- Other predictors (multiple labels per point), relationship with robust statistics, minimizing the number of queries
- Algorithms for (k, z) -clustering, i.e., $\sum_{p \in P} \min_{c \in C} \|p - c\|_2^z$
- Algorithms for L_p -metrics, i.e., $\sum_{p \in P} \min_{c \in C} \|p - c\|_p^p$