

Differentially Private k-Means with Constant Multiplicative Error

[Uri Stemmer](#)

Ben-Gurion University

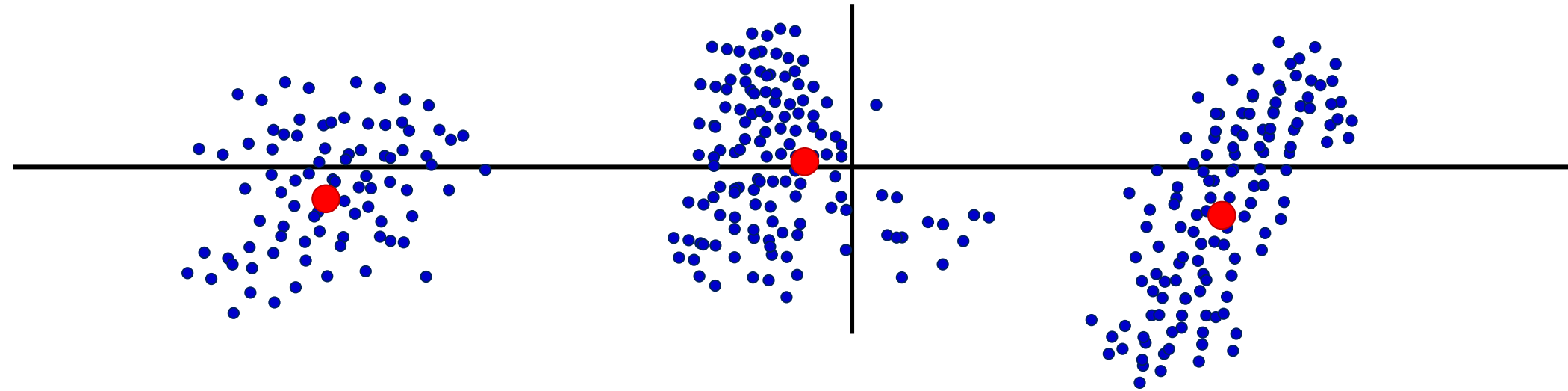
joint work with

Haim Kaplan

What is k -Means Clustering?

Given: Data points $S = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ and parameter k

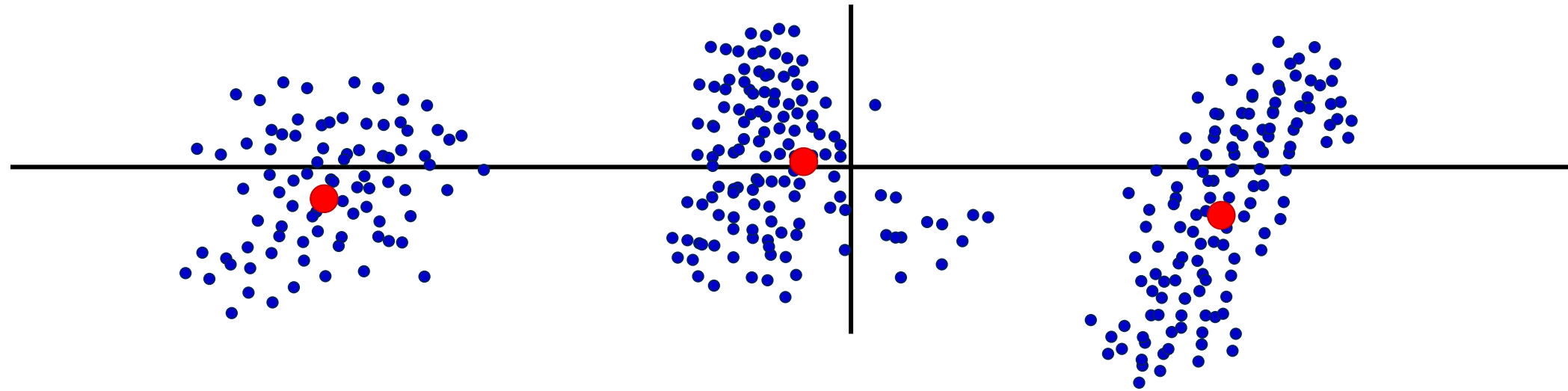
Identify k centers $C = (u_1, \dots, u_k)$ minimizing $\text{cost}(C) = \sum_i \min_\ell \|x_i - u_\ell\|^2$



What is k -Means Clustering?

Given: Data points $S = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ and parameter k

Identify k centers $C = (u_1, \dots, u_k)$ minimizing $\text{cost}(C) = \sum_i \min_\ell \|x_i - u_\ell\|^2$

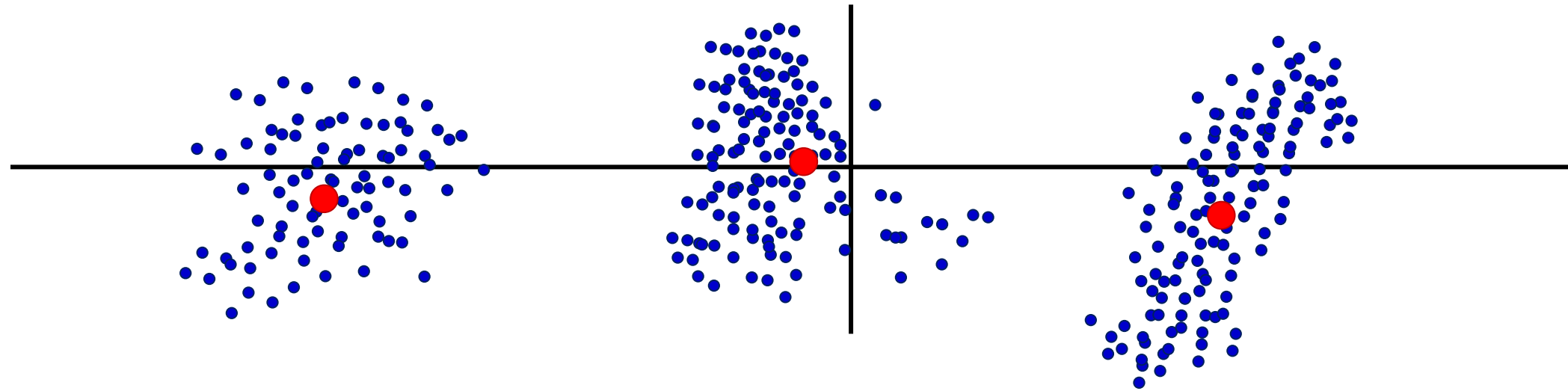


- ✓ Probably the most well-studied clustering problem
- ✓ Tons of applications
- ✓ Super popular

What is k -Means Clustering?

Given: Data points $S = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ and parameter k

Identify k centers $C = (u_1, \dots, u_k)$ minimizing $\text{cost}(C) = \sum_i \min_\ell \|x_i - u_\ell\|^2$



What is Differentially Private k -Means?

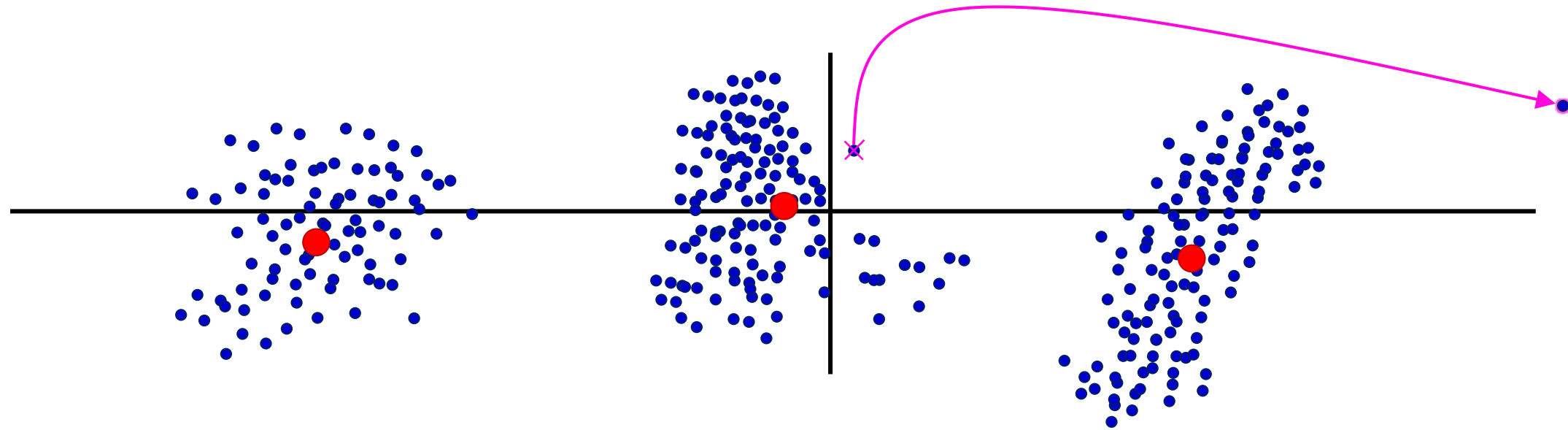
[Dwork, McSherry, Nissim, Smith 06] (informal)

- ✓ Every data point x_i represents the (private) information of one individual
- ✓ **Goal:** the output (the set of centers) does not reveal information that is specific to any single individual

What is k -Means Clustering?

Given: Data points $S = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ and parameter k

Identify k centers $C = (u_1, \dots, u_k)$ minimizing $\text{cost}(C) = \sum_i \min_\ell \|x_i - u_\ell\|^2$



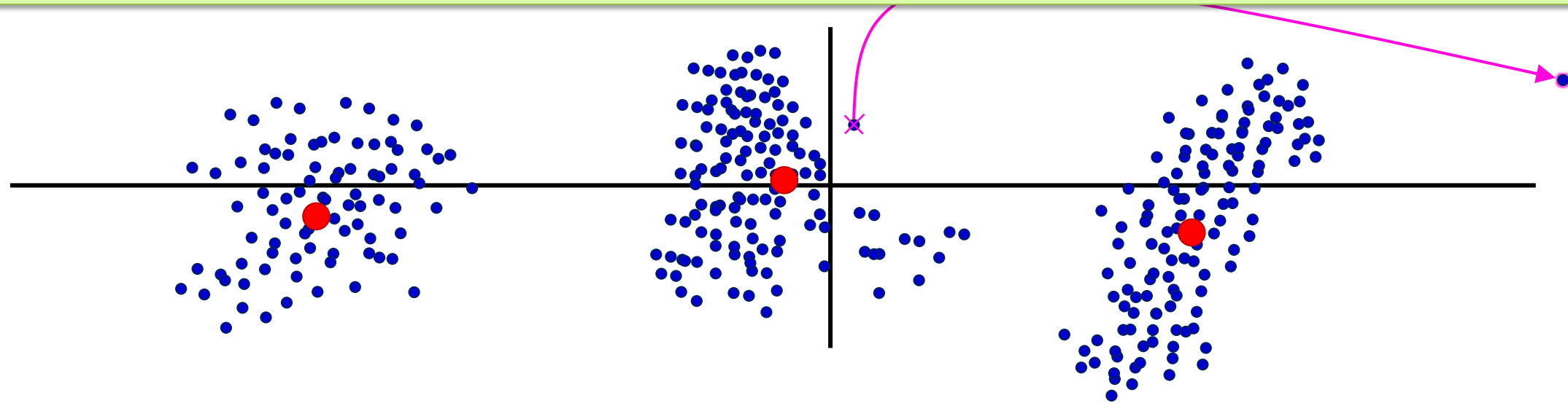
What is Differentially Private k -Means?

[Dwork, McSherry, Nissim, Smith 06] (informal)

- ✓ Every data point x_i represents the (private) information of one individual
- ✓ **Goal:** the output (the set of centers) does not reveal information that is specific to any single individual
- ✓ **Requirement:** the output distribution is insensitive to any arbitrarily change of a single input point (an algorithm satisfying this requirement is *differentially private*)

Why is that a good privacy definition?

Even if an observer knows all other data point but mine, and now she sees the outcome of the computation, then she still cannot learn “anything” on my data point



What is Differentially Private *k*-Means?

[Dwork, McSherry, Nissim, Smith 06] (informal)

- ✓ Every data point x_i represents the (private) information of one individual
- ✓ **Goal:** the output (the set of centers) does not reveal information that is specific to any single individual
- ✓ **Requirement:** the output distribution is insensitive to any arbitrarily change of a single input point (an algorithm satisfying this requirement is *differentially private*)

Differentially Private k -Means Clustering

Given: Data points $\mathcal{S} = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ and parameter k

Identify k centers $\mathcal{C} = (u_1, \dots, u_k)$ minimizing $\text{cost}(\mathcal{C}) = \sum_i \min_{\ell} \|x_i - u_{\ell}\|^2$

Requirement: the output distribution is insensitive to any arbitrarily change of a single input point

Differentially Private k -Means Clustering

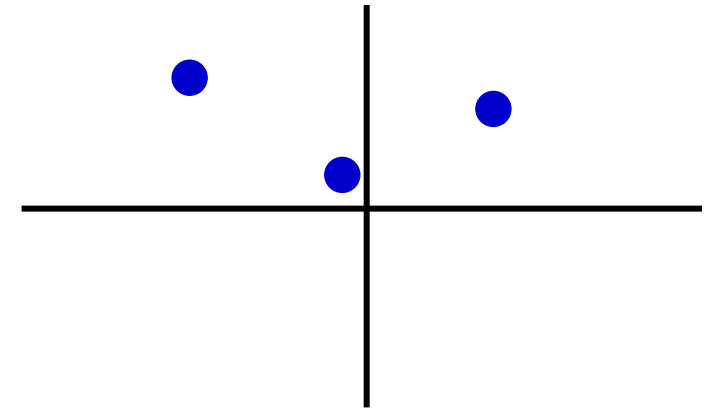
Given: Data points $\mathcal{S} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in (\mathbb{R}^d)^n$ and parameter k

Identify k centers $\mathcal{C} = (\mathbf{u}_1, \dots, \mathbf{u}_k)$ minimizing $\text{cost}(\mathcal{C}) = \sum_i \min_{\ell} \|\mathbf{x}_i - \mathbf{u}_{\ell}\|^2$

Requirement: the output distribution is insensitive to any arbitrarily change of a single input point

Observe: With privacy we must have additive error

- Assume $k = n = 3$
- OPT's cost = 0



Differentially Private k -Means Clustering

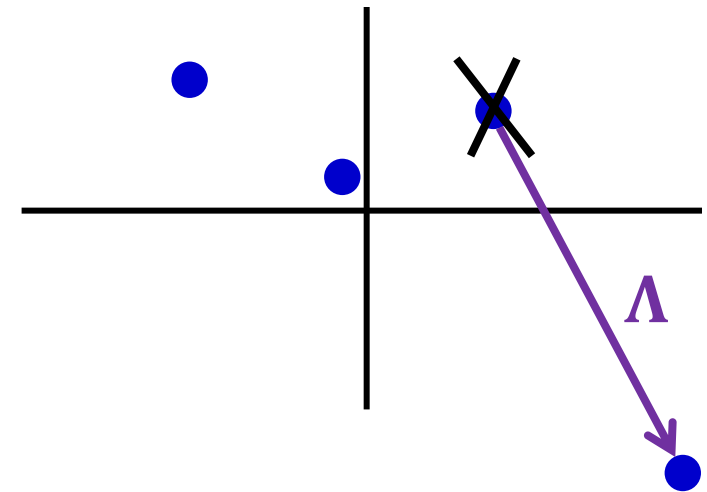
Given: Data points $S = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ and parameter k

Identify k centers $C = (u_1, \dots, u_k)$ minimizing $\text{cost}(C) = \sum_i \min_\ell \|x_i - u_\ell\|^2$

Requirement: the output distribution is insensitive to any arbitrarily change of a single input point

Observe: With privacy we must have additive error

- Assume $k = n = 3$
- OPT's cost = 0
- Move one point
- OPT's cost = 0



Differentially Private k -Means Clustering

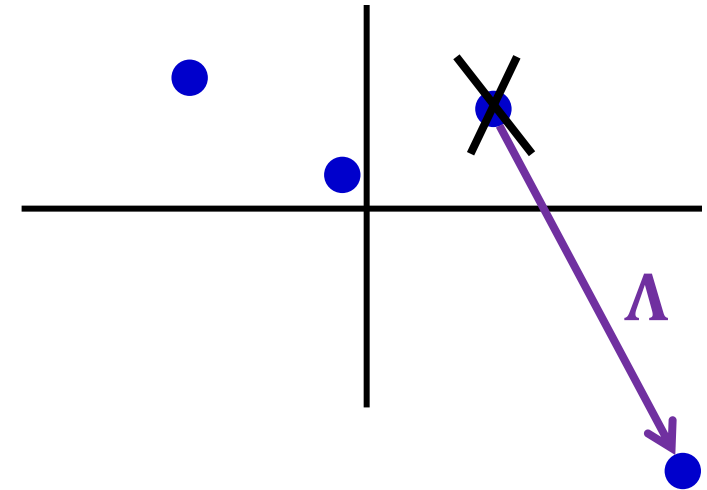
Given: Data points $S = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ and parameter k

Identify k centers $C = (u_1, \dots, u_k)$ minimizing $\text{cost}(C) = \sum_i \min_\ell \|x_i - u_\ell\|^2$

Requirement: the output distribution is insensitive to any arbitrarily change of a single input point

Observe: With privacy we must have additive error

- Assume $k = n = 3$
- OPT's cost = 0
- Move one point
- OPT's cost = 0
- Each solution must remain approx. equally likely
- On at least one of these inputs our cost is $\approx \Lambda^2$



Differentially Private k -Means Clustering

Given: Data points $S = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ and parameter k

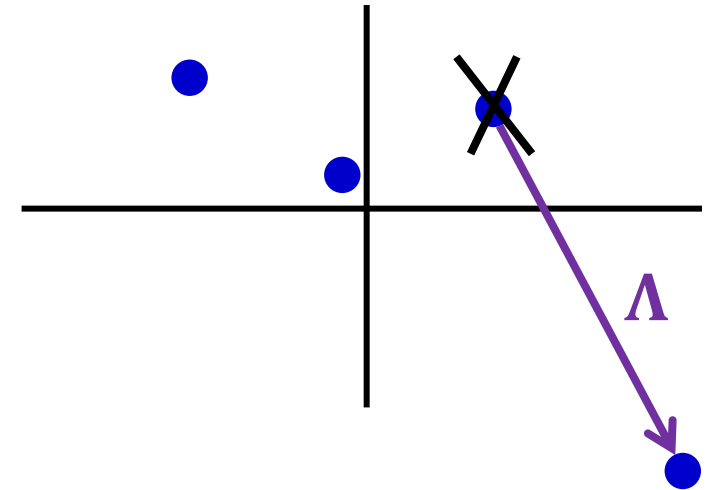
Identify k centers $C = (u_1, \dots, u_k)$ minimizing $\text{cost}(C) = \sum_i \min_\ell \|x_i - u_\ell\|^2$

Requirement: the output distribution is insensitive to any arbitrarily change of a single input point

Observe: With privacy we must have additive error

- Assume $k = n = 3$
- OPT's cost = 0
- Move one point
- OPT's cost = 0
- Each solution must remain approx. equally likely
- On at least one of these inputs our cost is $\approx \Lambda^2$

\Rightarrow We assume that input points come from the unit ball



Previous and New Bounds

Ref	Model	Runtime	Bounds
GLMRT'10	differential privacy	n^d	$O(1) \cdot \text{OPT} + \tilde{O}(k^2 \cdot d)$
NCBN'16	differential privacy	poly	$O(\log k) \cdot \text{OPT} + \tilde{O}(n)$
FXZR'17	differential privacy	poly	$O(k \log n) \cdot \text{OPT} + \tilde{O}(k^{3/2} \cdot \sqrt{d})$
BDLMZ'17	differential privacy	poly	$O(\log^3 n) \cdot \text{OPT} + \tilde{O}(k^2 + d)$
NS'18	differential privacy	poly	$O(k) \cdot \text{OPT} + \tilde{O}(k^{1.51} \cdot d^{0.51})$
New	differential privacy	poly	$O(1) \cdot \text{OPT} + \tilde{O}(k^{1.01} \cdot d^{0.51} + k^{3/2})$