# 6.036/6.862: Introduction to Machine Learning

**Lecture:** starts Tuesdays 9:35am (Boston time zone)

**Course website:** introml.odl.mit.edu

**Who's talking?** Prof. Tamara Broderick

**Questions?** discourse.odl.mit.edu ("Lecture 13" category)

**Materials:** Will all be available at course website

**Last Time(s)**

I.  Supervised Learning

- Classification

- Regression

**Today's Plan**

I.   Unsupervised learning

II.  Clustering

III. k-means clustering

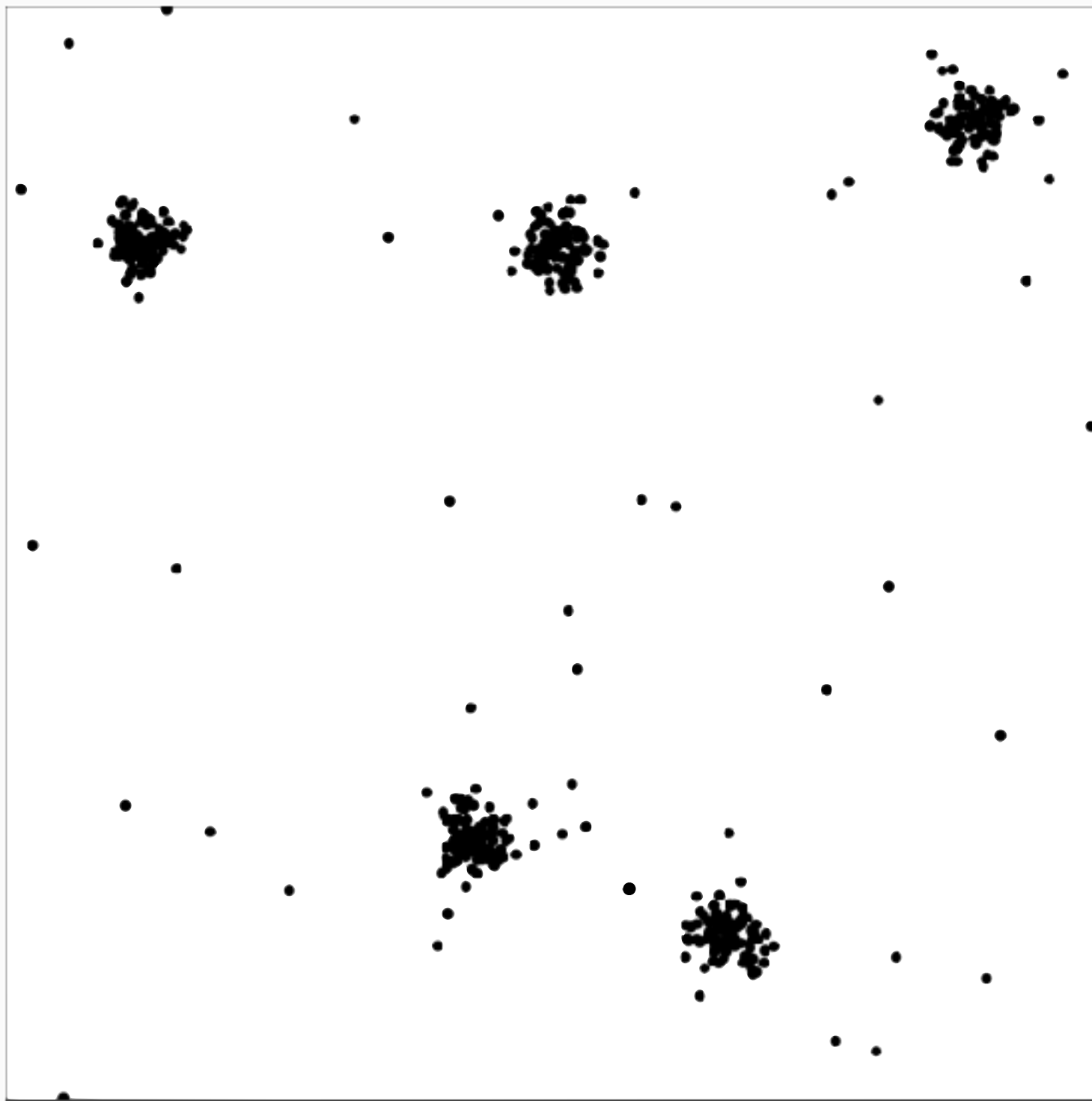# Food distribution placement

# Food distribution placement



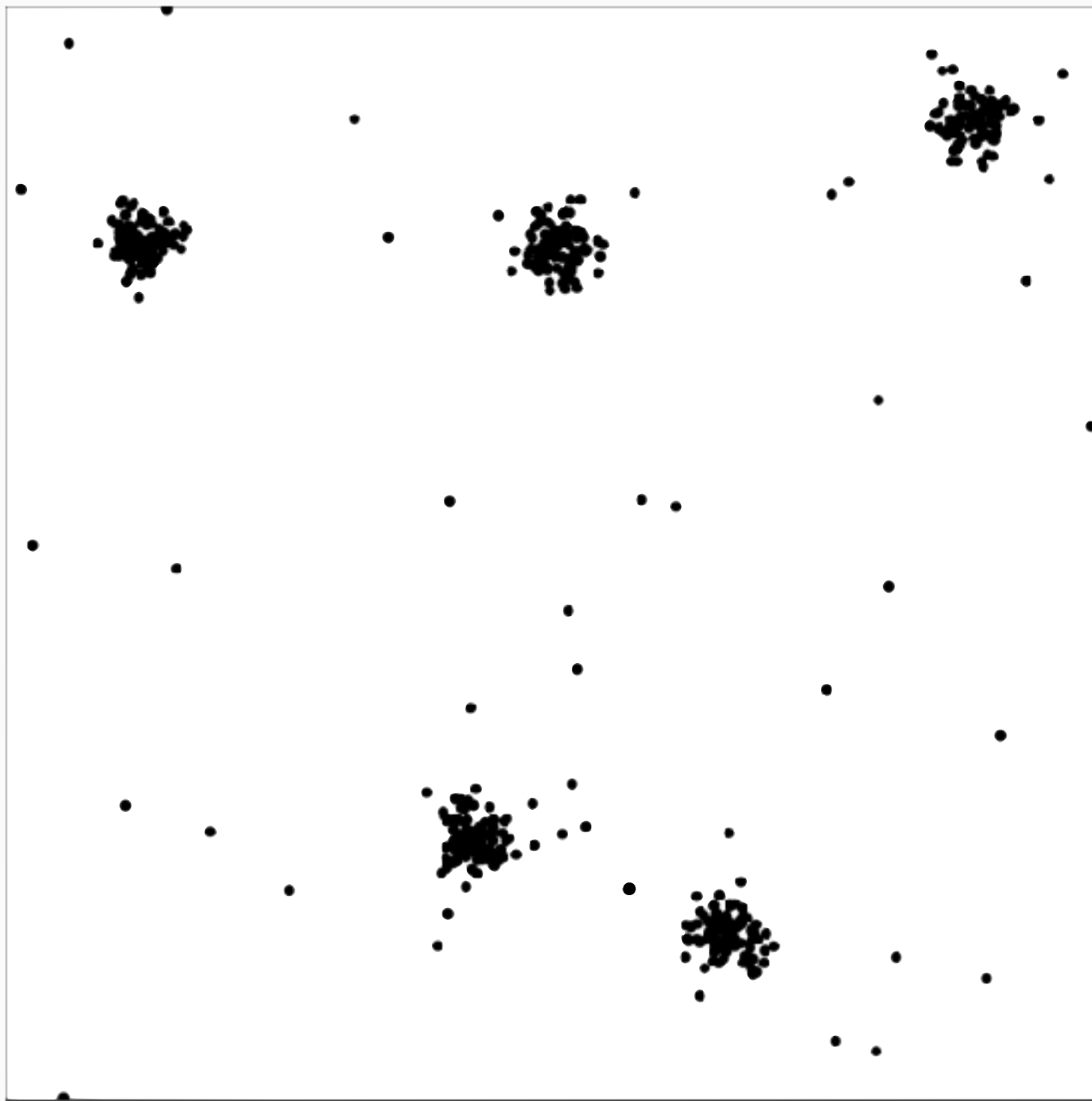Yes Free Lunch

# Food distribution placement

# Food distribution placement

- Where should I have my $k$ food trucks park?

# Food distribution placement
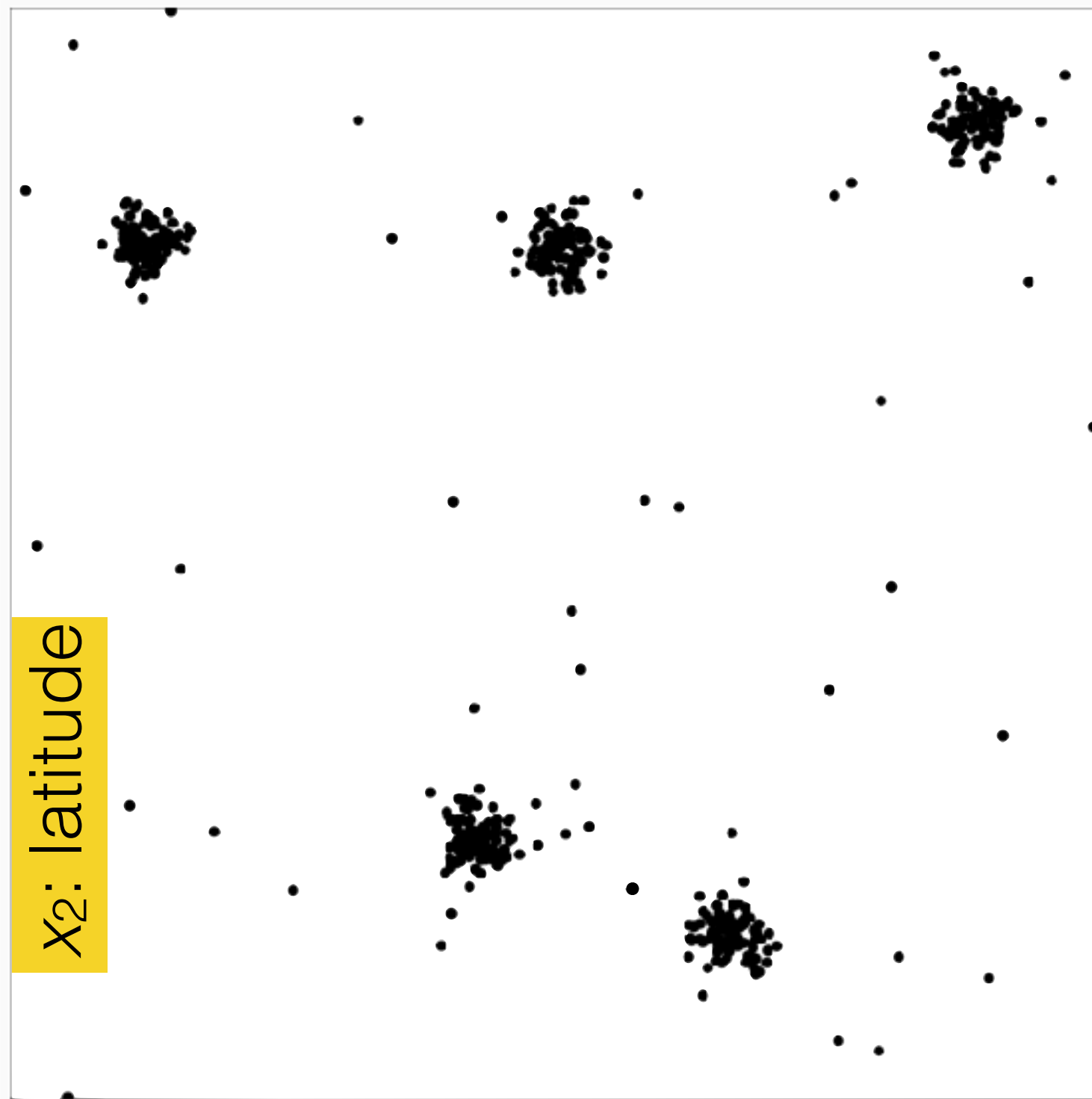


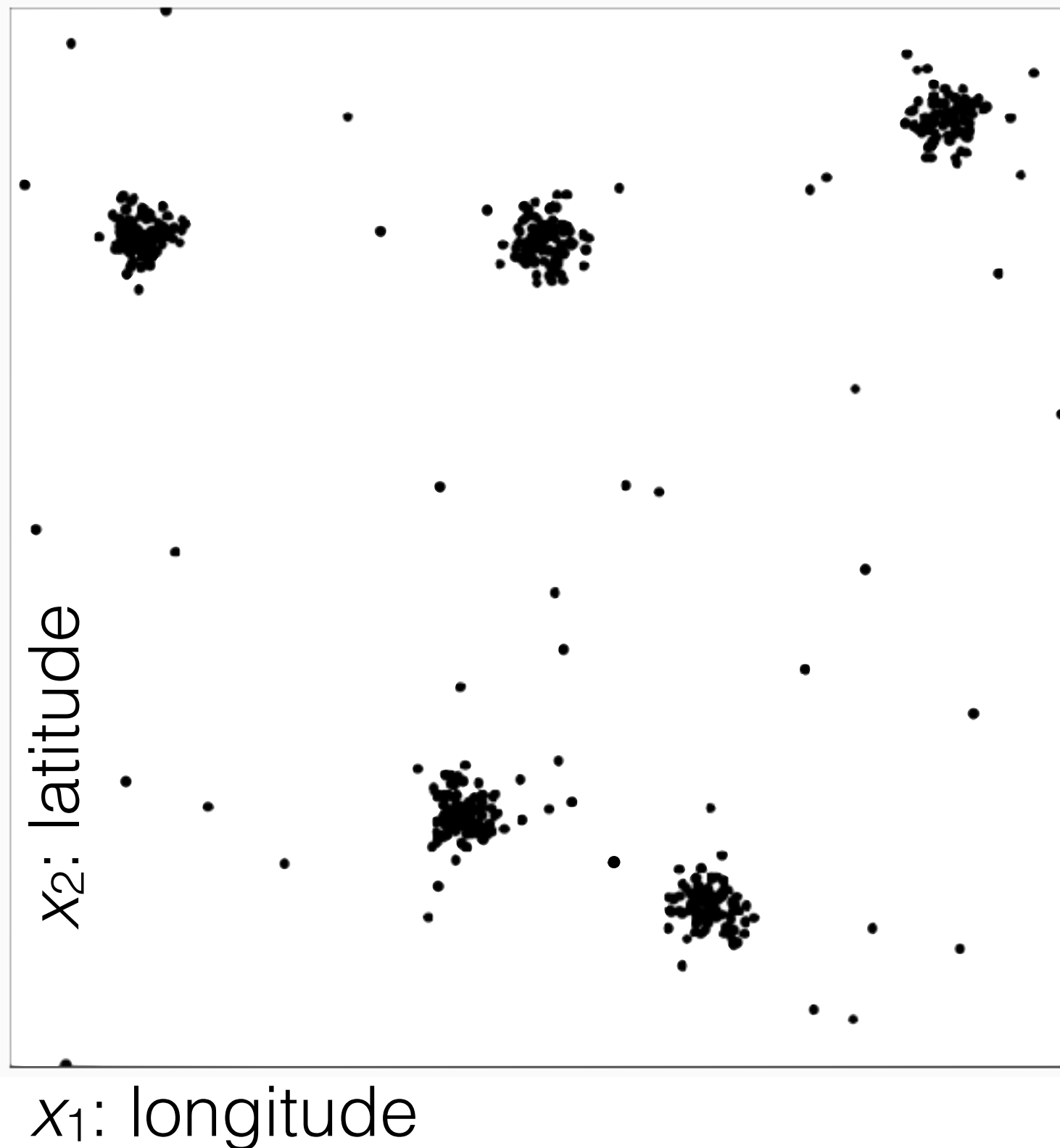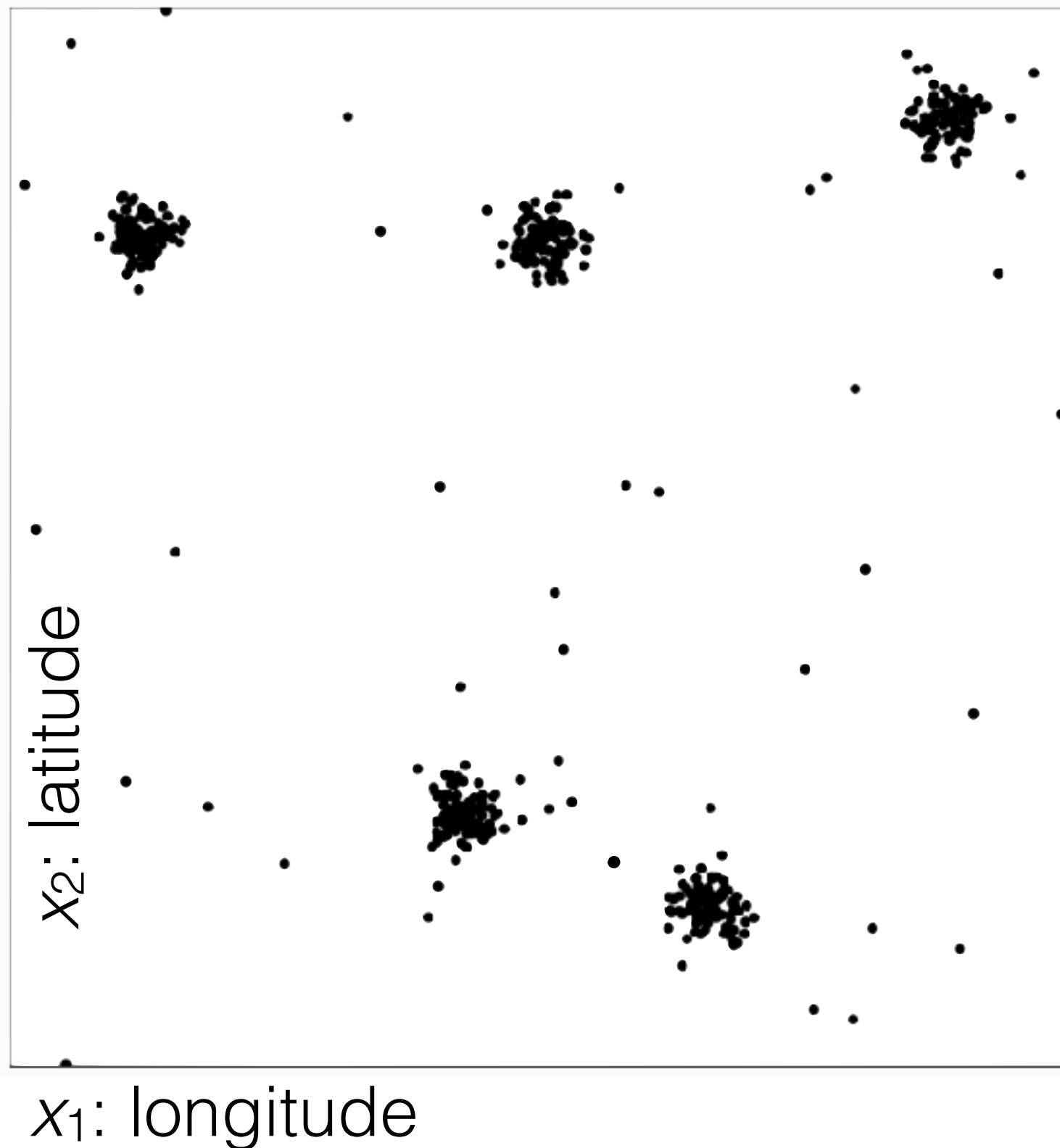- Where should I have my $k$ food trucks park?

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]

# Food distribution placement



$x_1$: longitude

- Where should I have my $k$ food trucks park?

# Food distribution placement



$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]
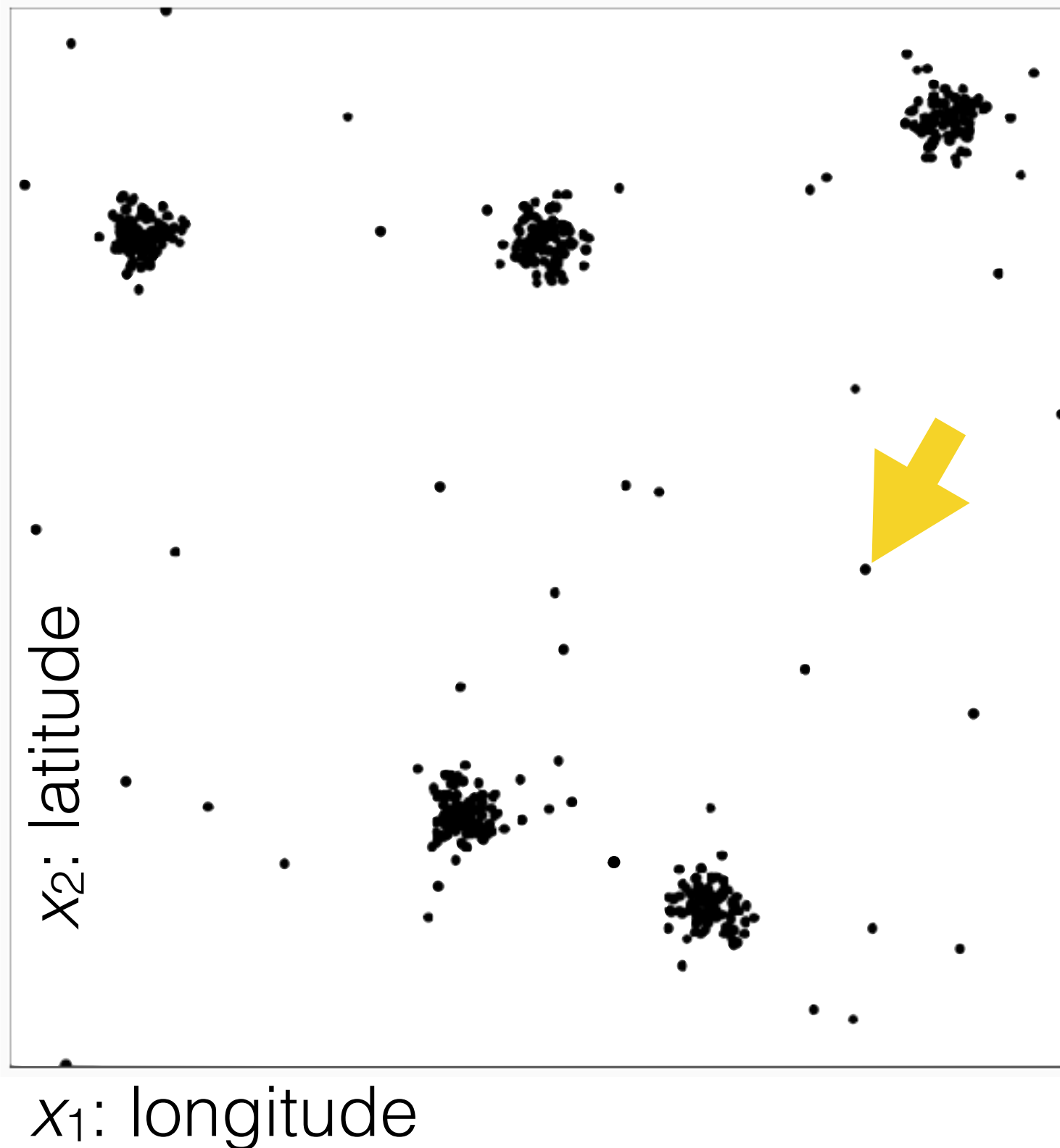
# Food distribution placement



$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?

- Want to minimize the loss of people we serve

# Food distribution placement



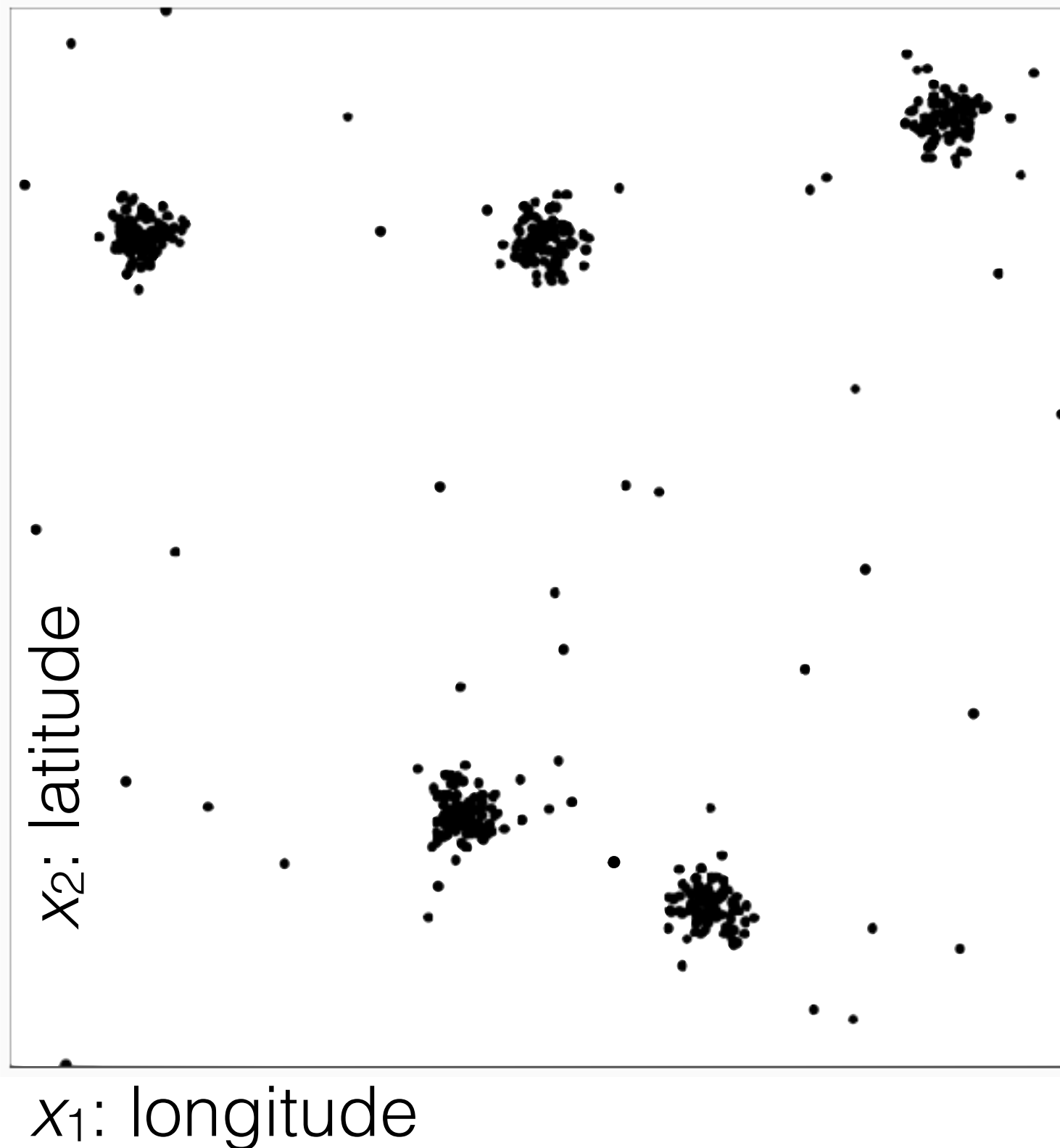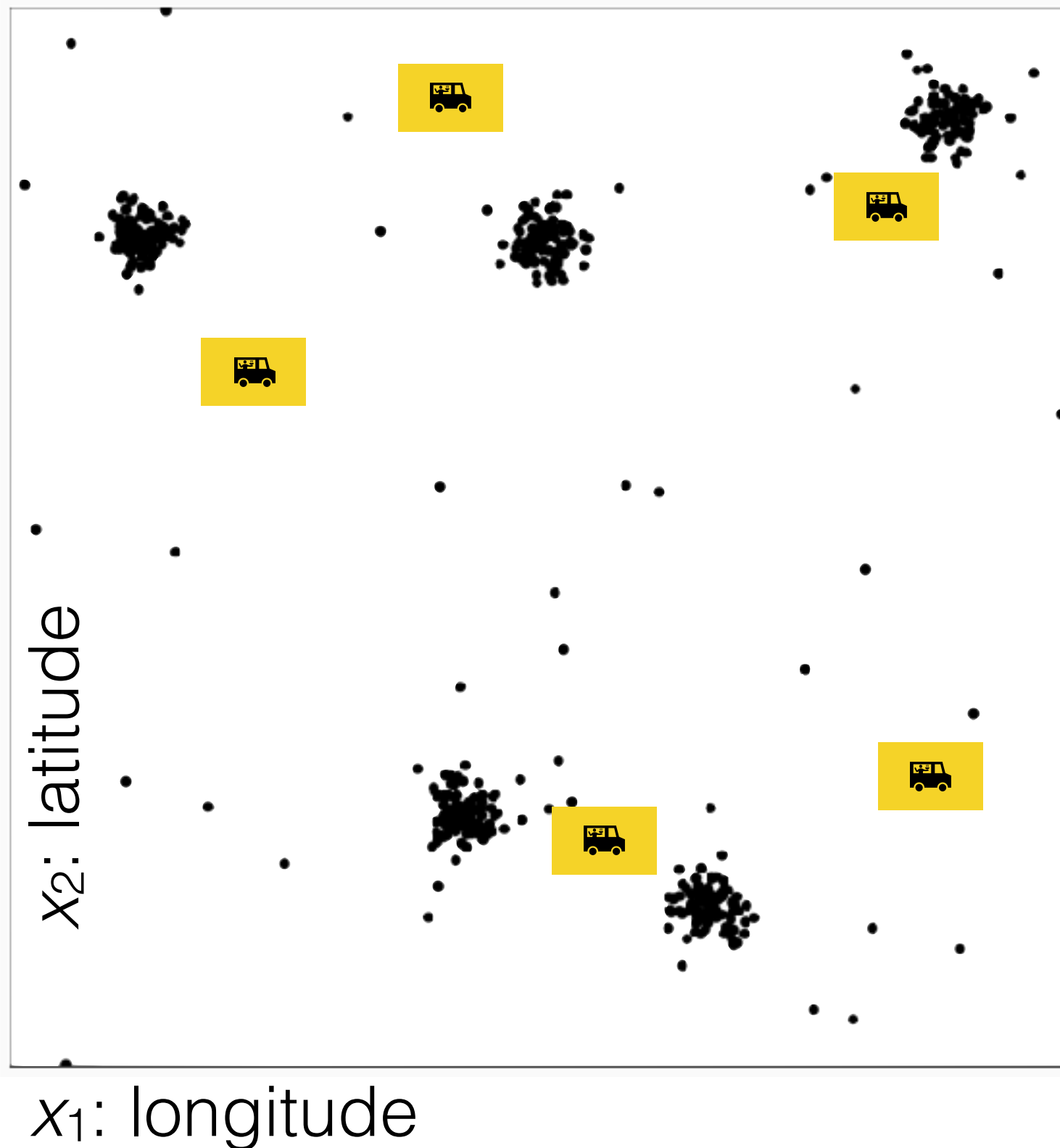$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$

3

# Food distribution placement



- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$

$x_2$: latitude

$x_1$: longitude

3

# Food distribution placement



$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]
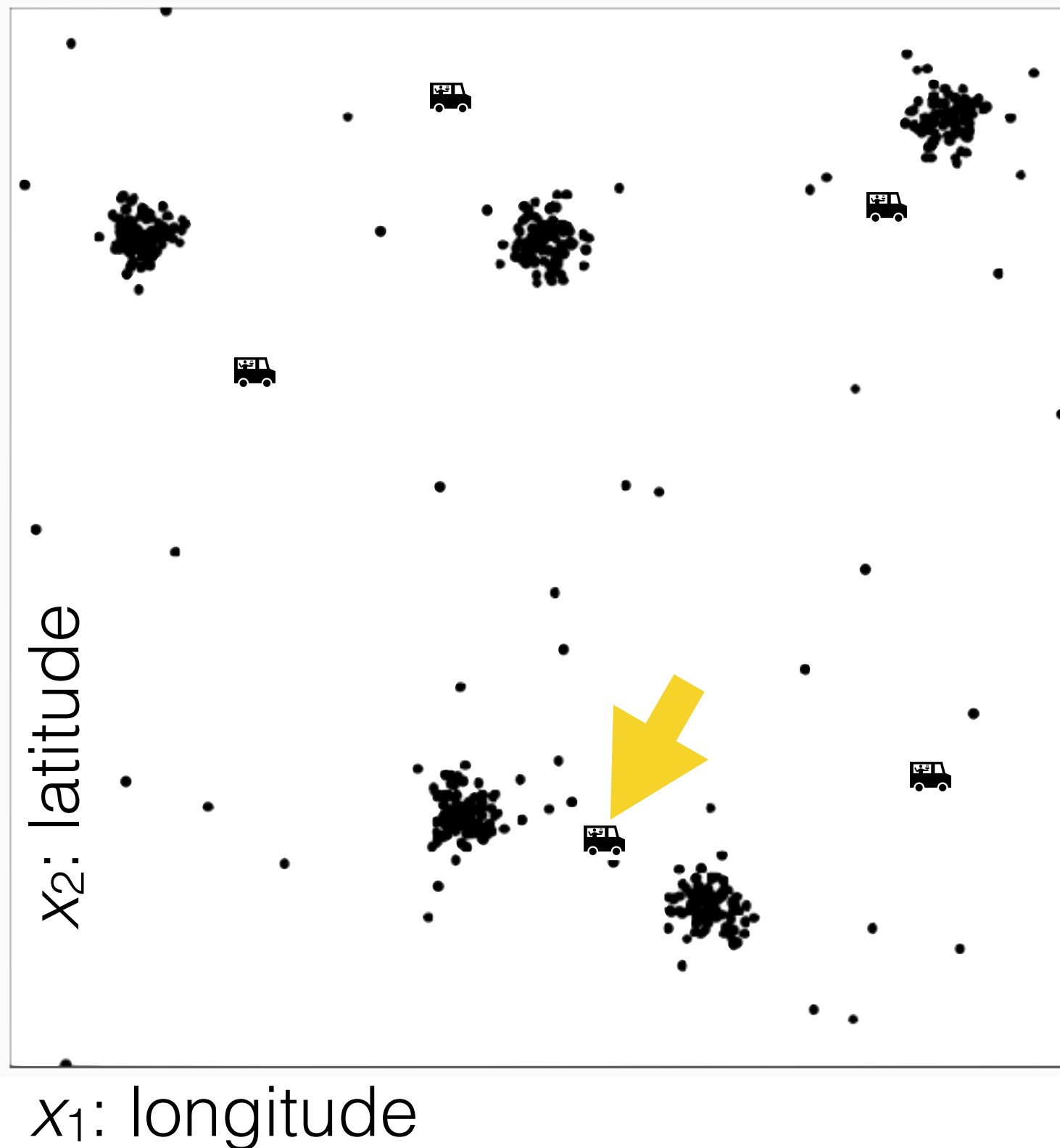
# Food distribution placement



- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
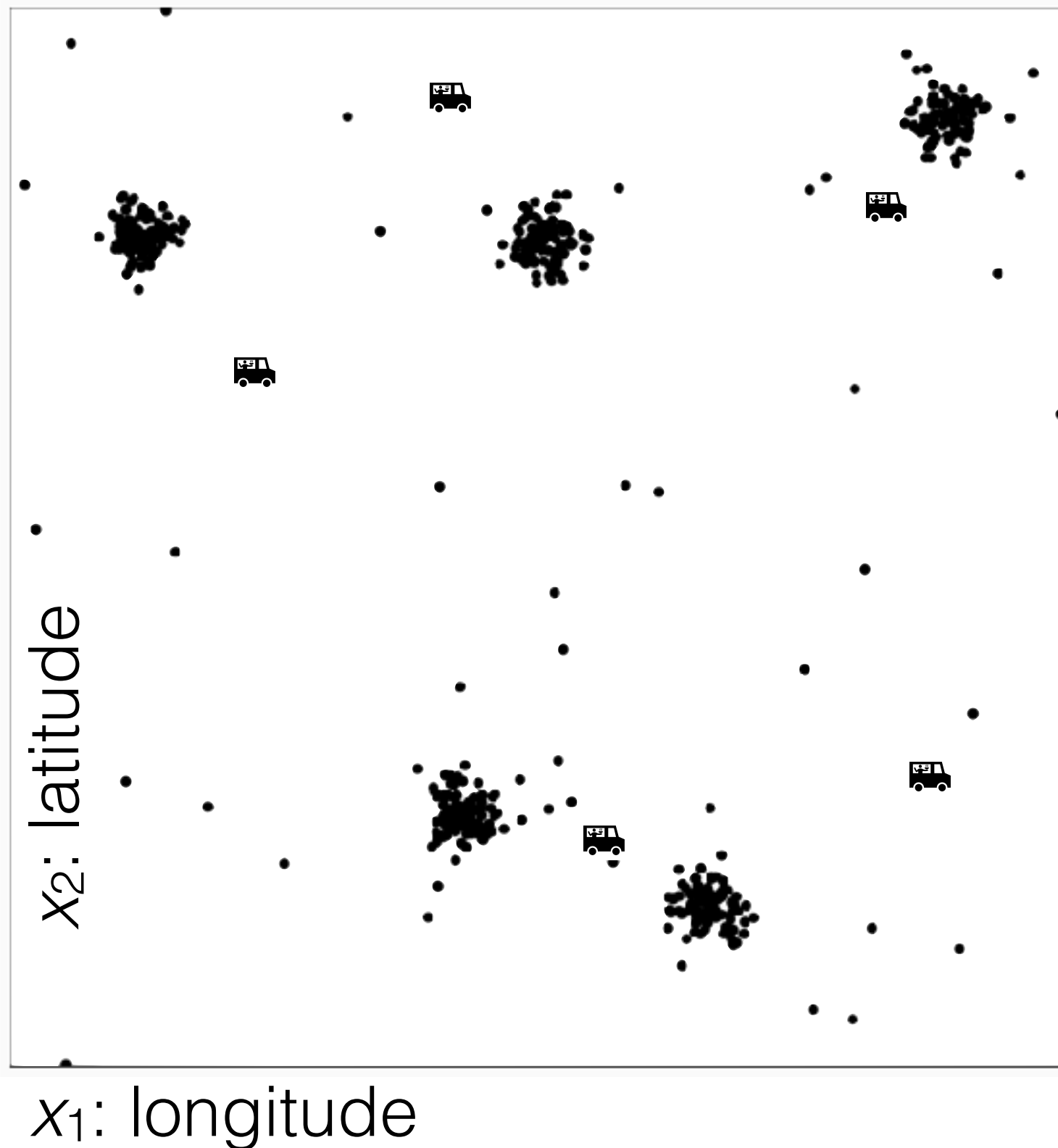
3

# Food distribution placement



- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$

$x_2$: latitude

$x_1$: longitude

3

# Food distribution placement
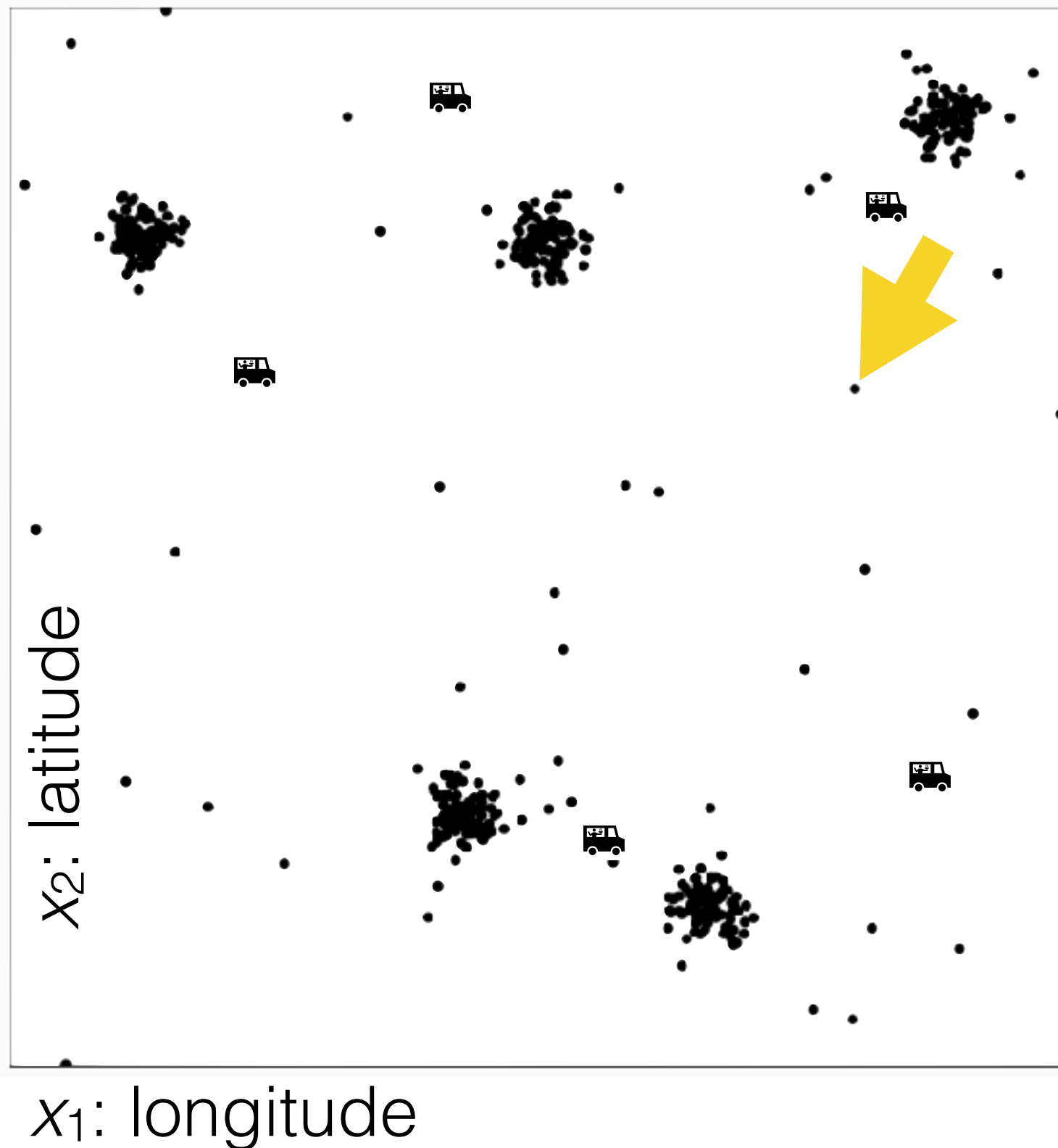


$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$

3

# Food distribution placement



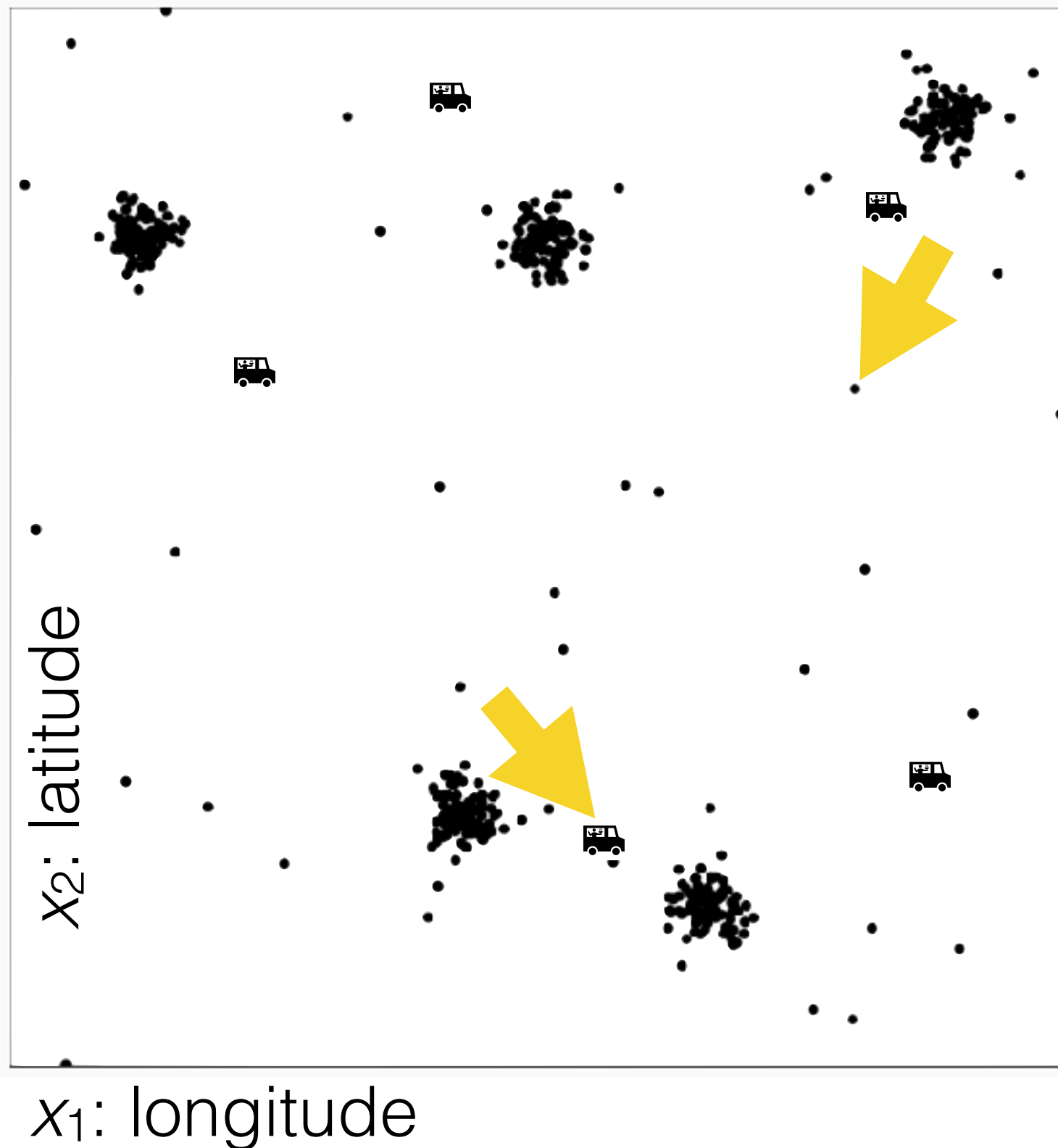$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$

3

# Food distribution placement



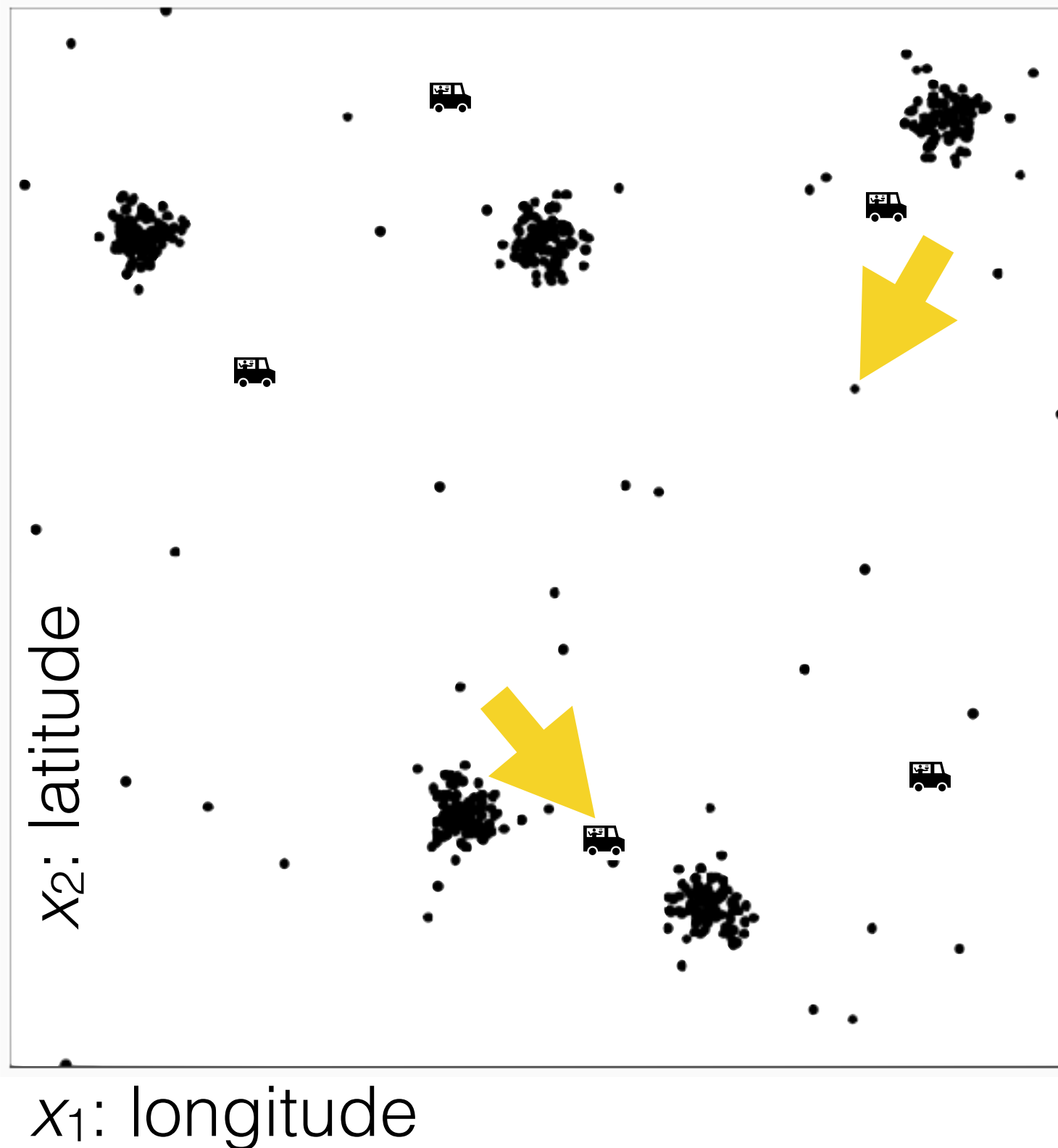$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$

3

# Food distribution placement



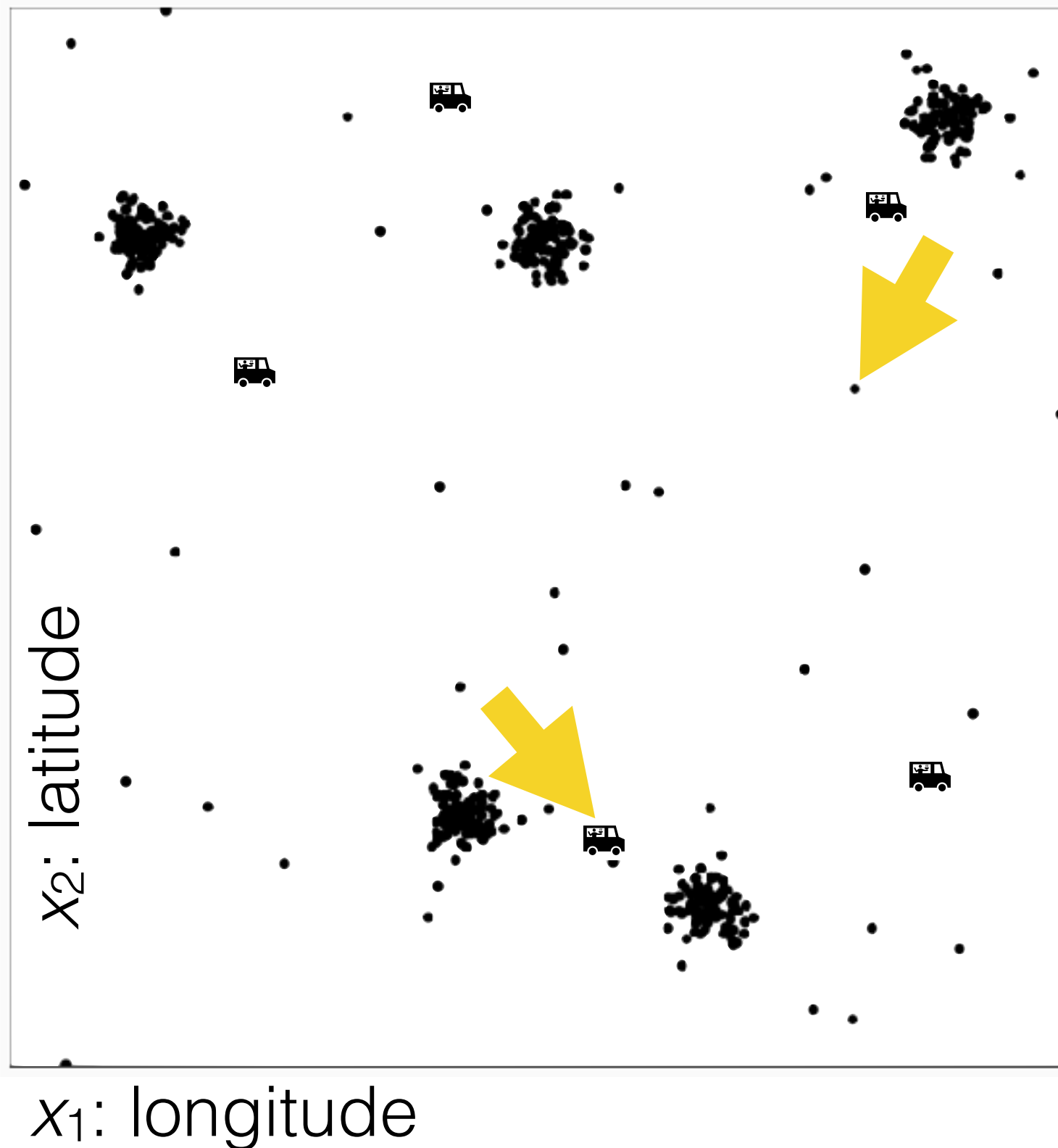$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
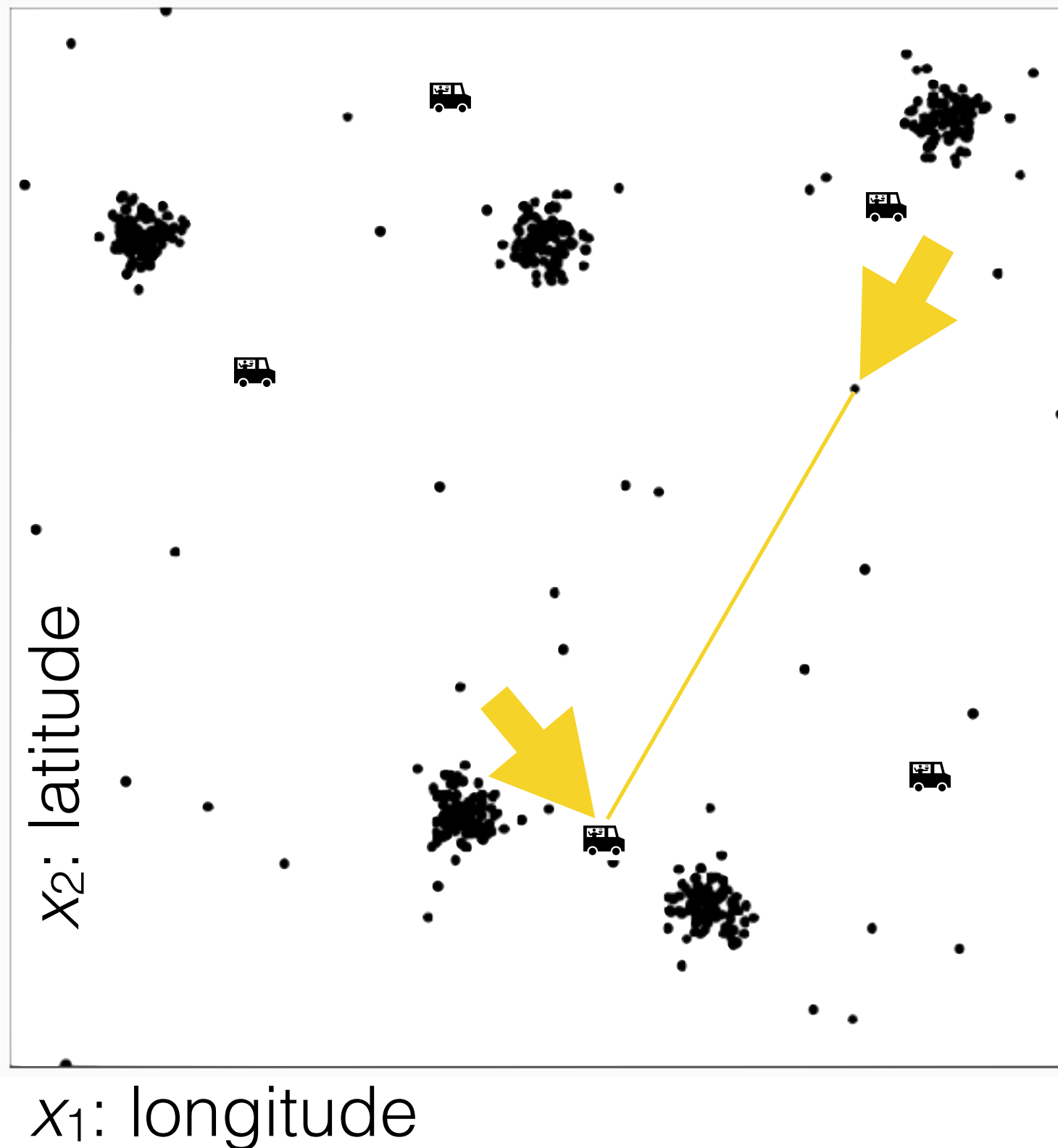- Loss if $i$ walks to truck $j$:

# Food distribution placement



- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:

$$\|x^{(i)} - \mu^{(j)}\|_2^2$$
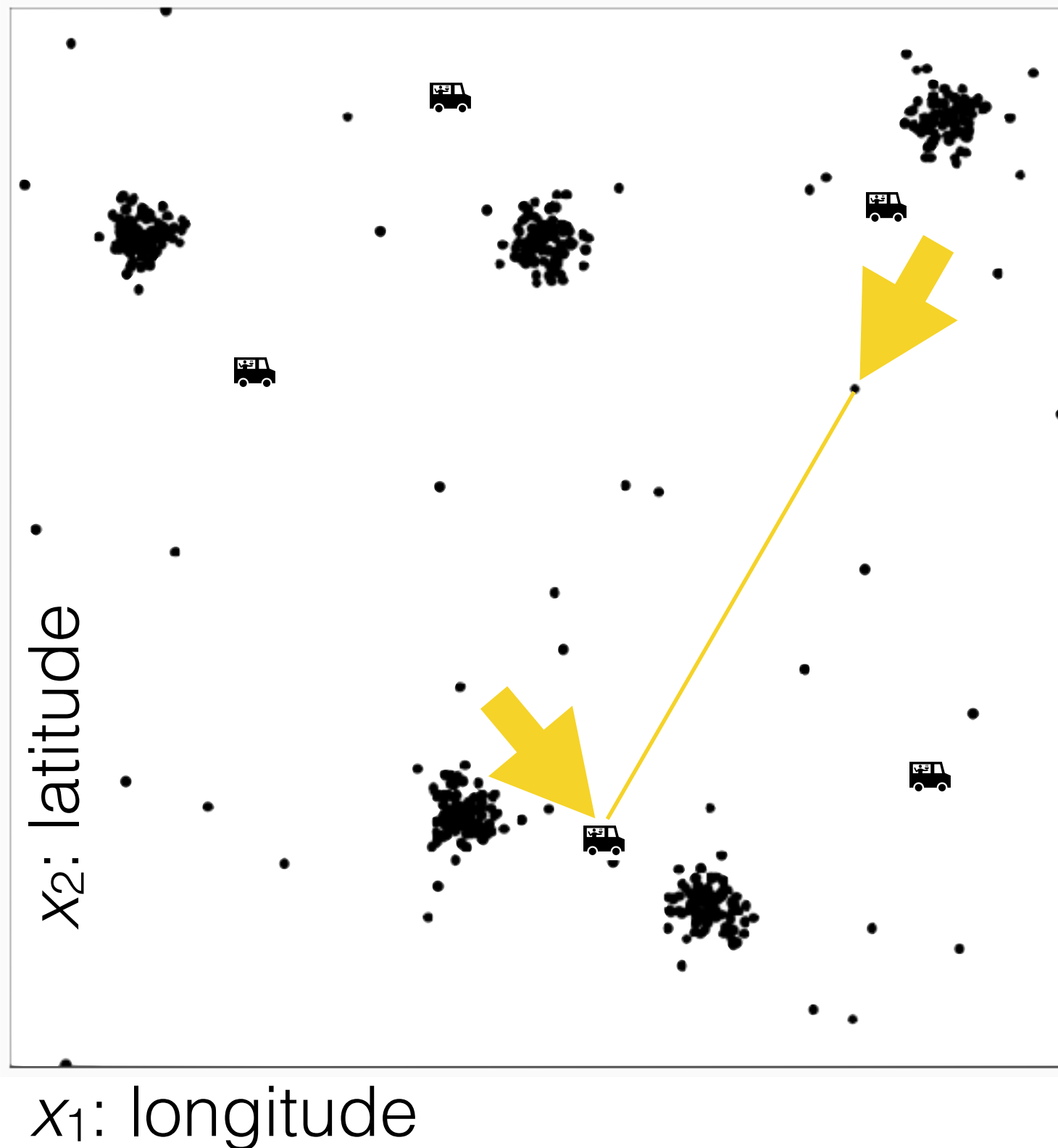
$x_2$: latitude

$x_1$: longitude

3

# Food distribution placement



- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:

$$\|x^{(i)} - \mu^{(j)}\|_2^2$$
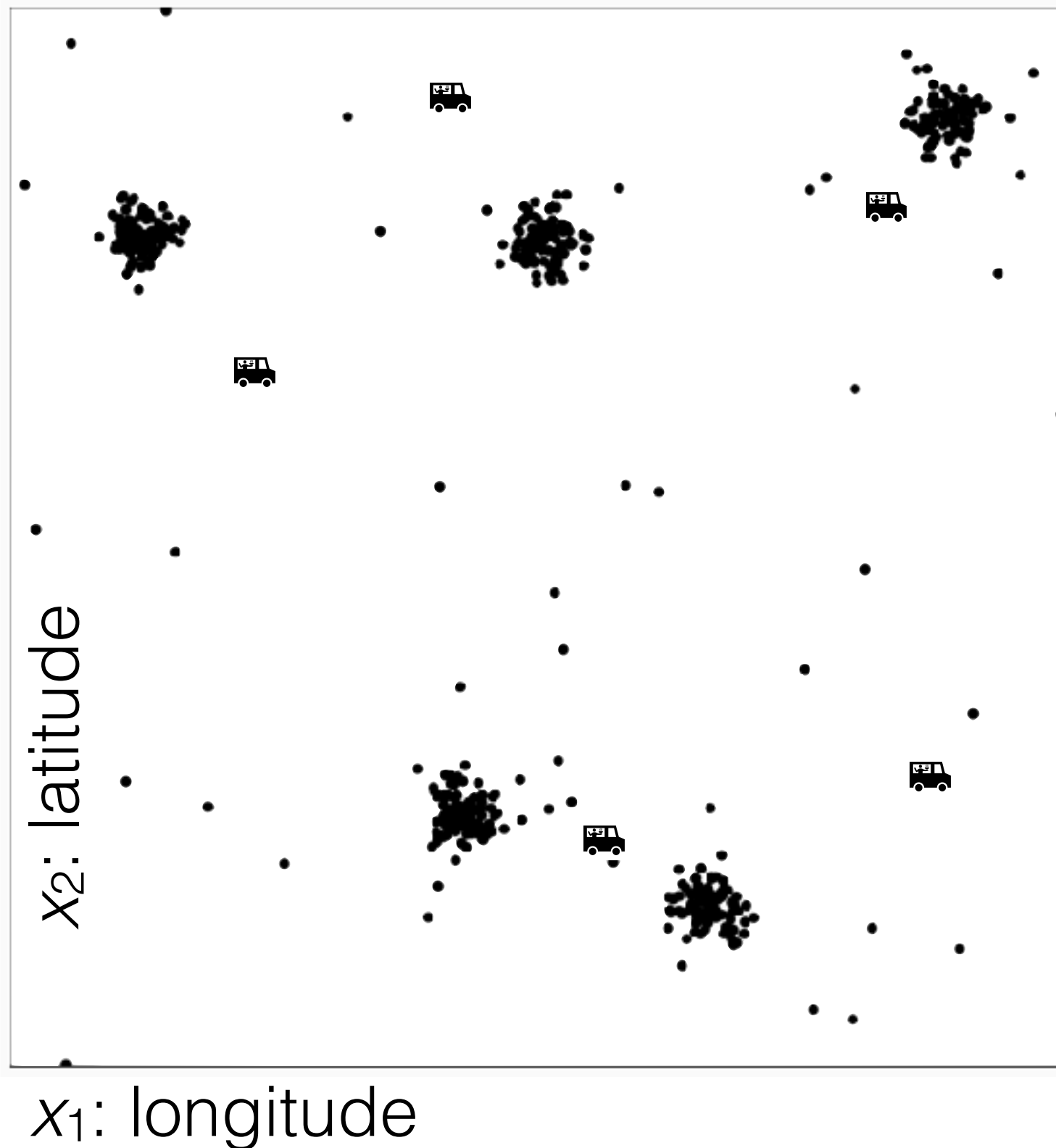
$x_2$: latitude

$x_1$: longitude

3

# Food distribution placement



- Where should I have my *k* food trucks park?
- Want to minimize the loss of people we serve
- Person *i* location $x^{(i)}$
- Food truck *j* location $\mu^{(j)}$
- Index of truck where person *i* walks: $y^{(i)}$
- Loss if *i* walks to truck *j*:

$$\|x^{(i)} - \mu^{(j)}\|_2^2$$

$x_1$: longitude

$x_2$: latitude

3

# Food distribution placement



$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
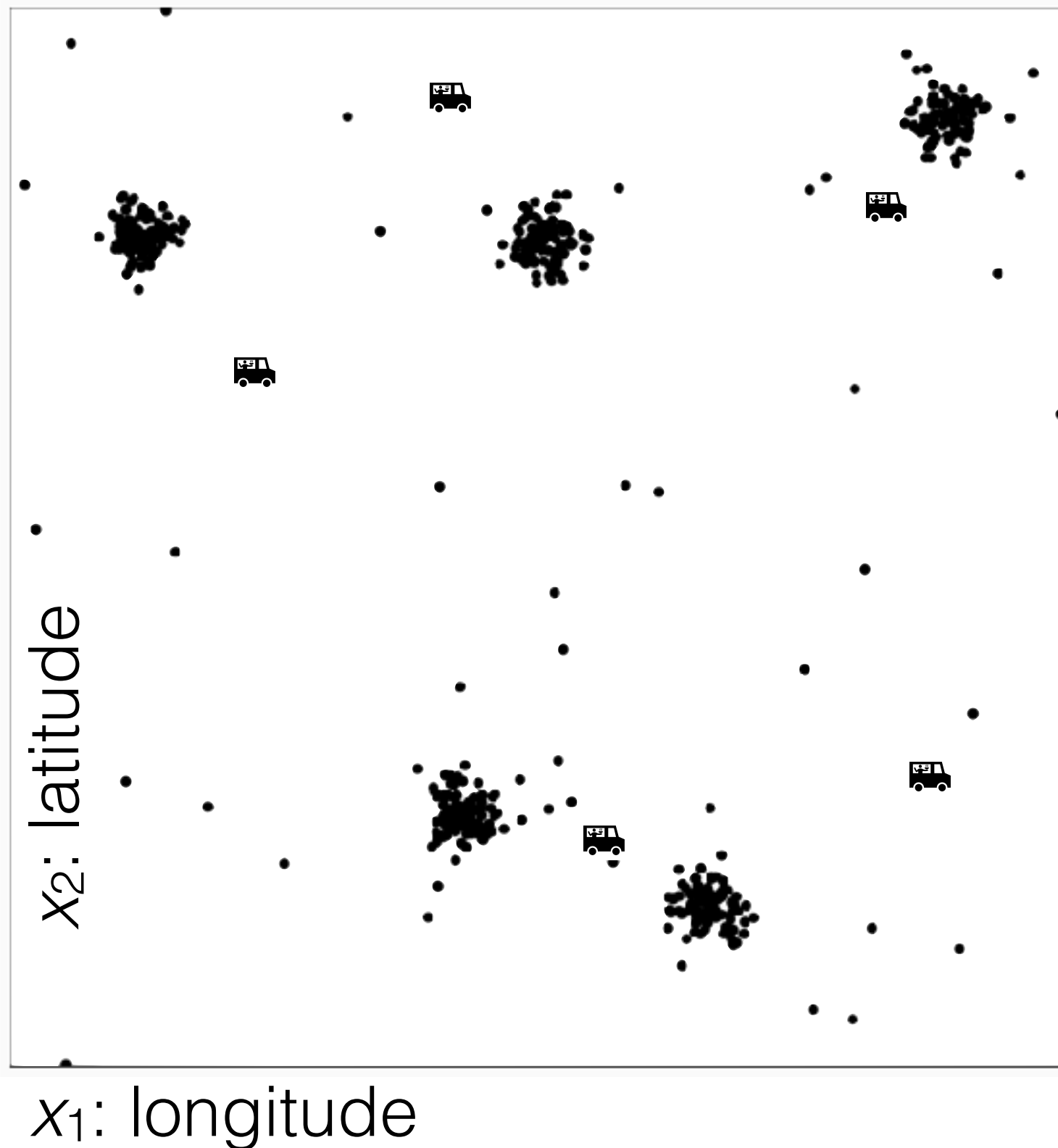- Loss if $i$ walks to truck $j$:
$$\|x^{(i)} - \mu^{(j)}\|_2^2$$
- Loss across all people:

3

# Food distribution placement
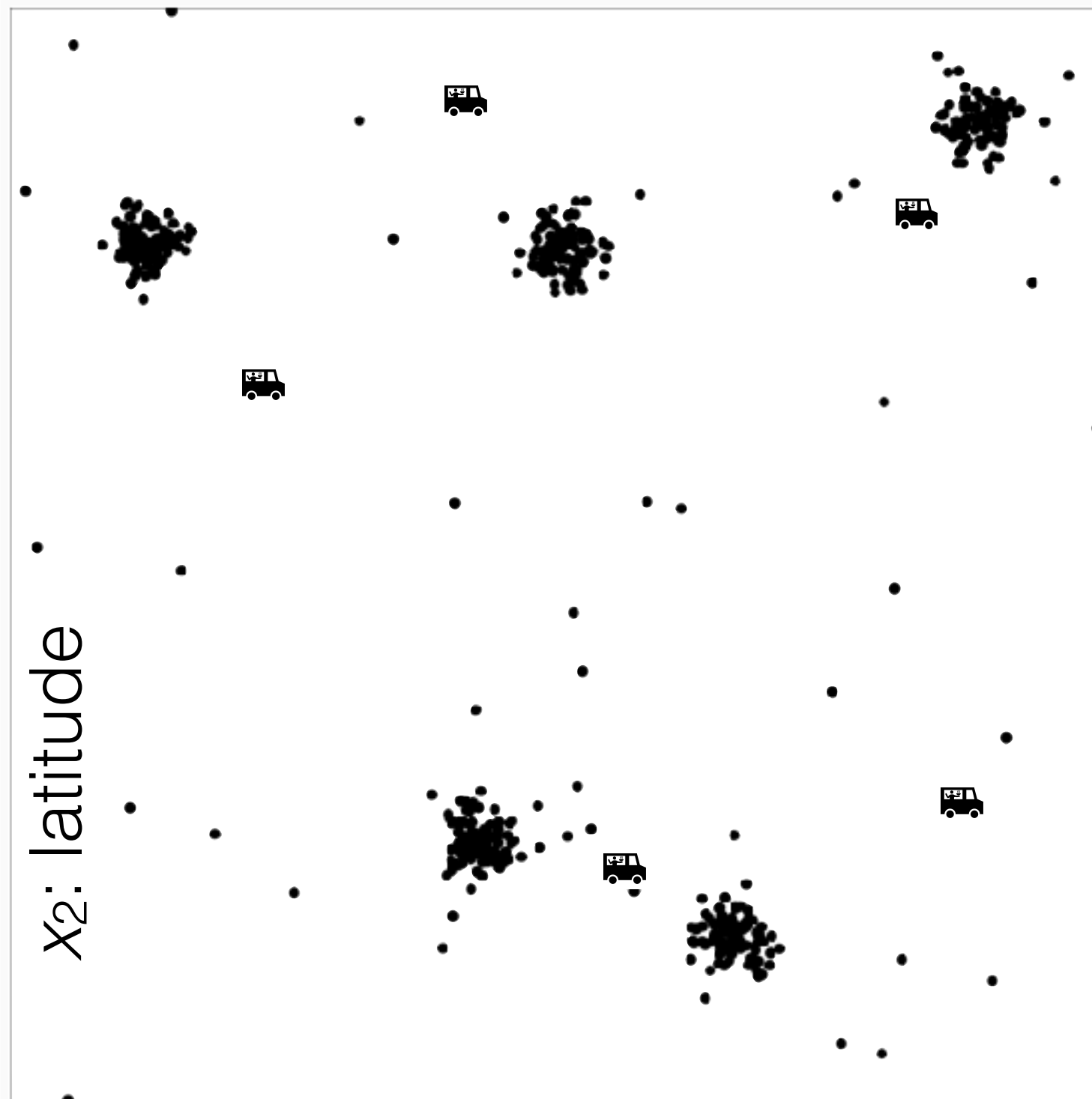


$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:

$$\|x^{(i)} - \mu^{(j)}\|_2^2$$

- Loss across all people:

$$\sum_{i=1}^{n} \|x^{(i)} - \mu^{(y^{(i)})}\|_2^2$$

3

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]

# Food distribution placement
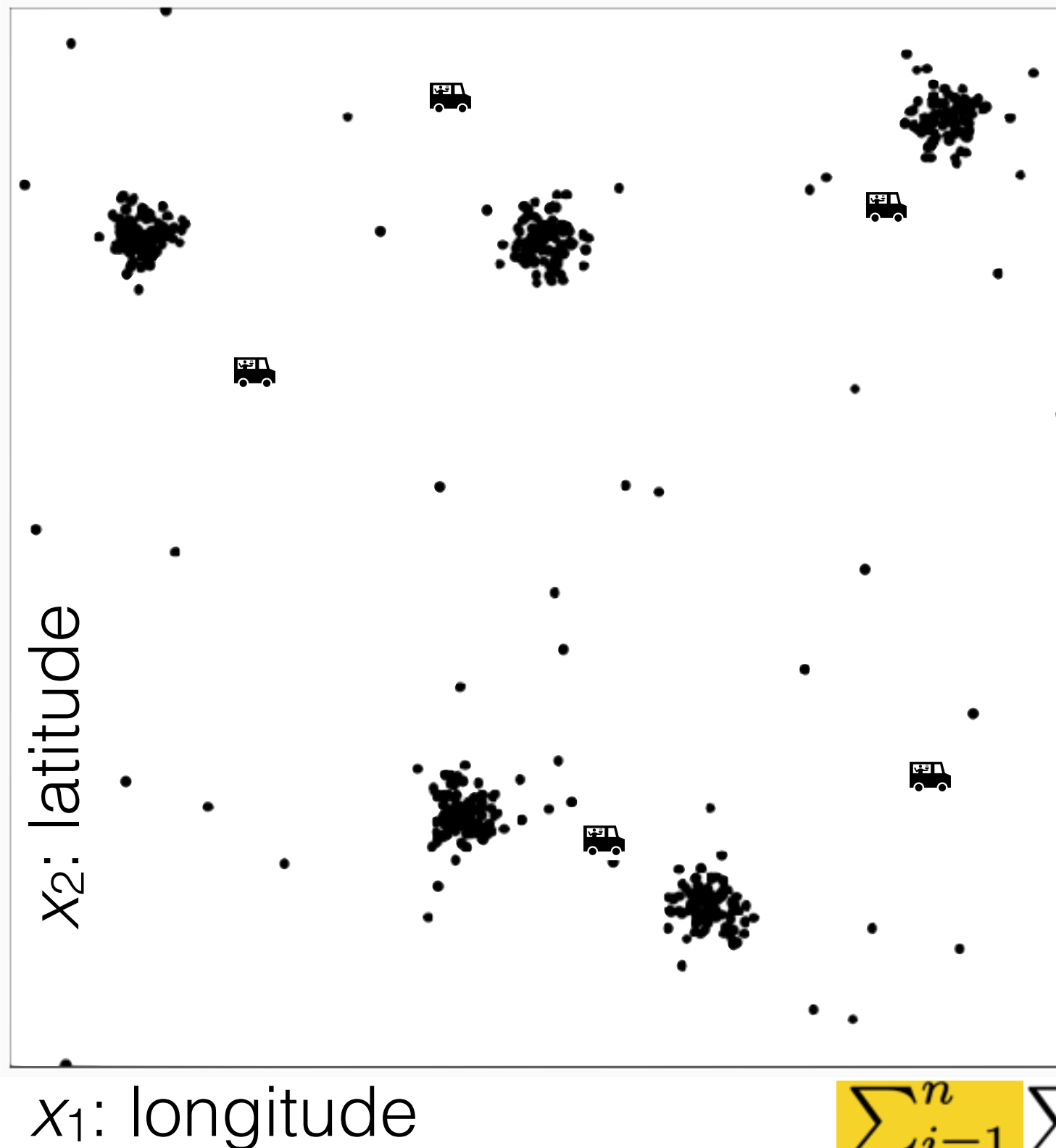


$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:
$$\|x^{(i)} - \mu^{(j)}\|_2^2$$
- Loss across all people:
$$\sum_{i=1}^{n} \sum_{j=1}^{k} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

3

# Food distribution placement



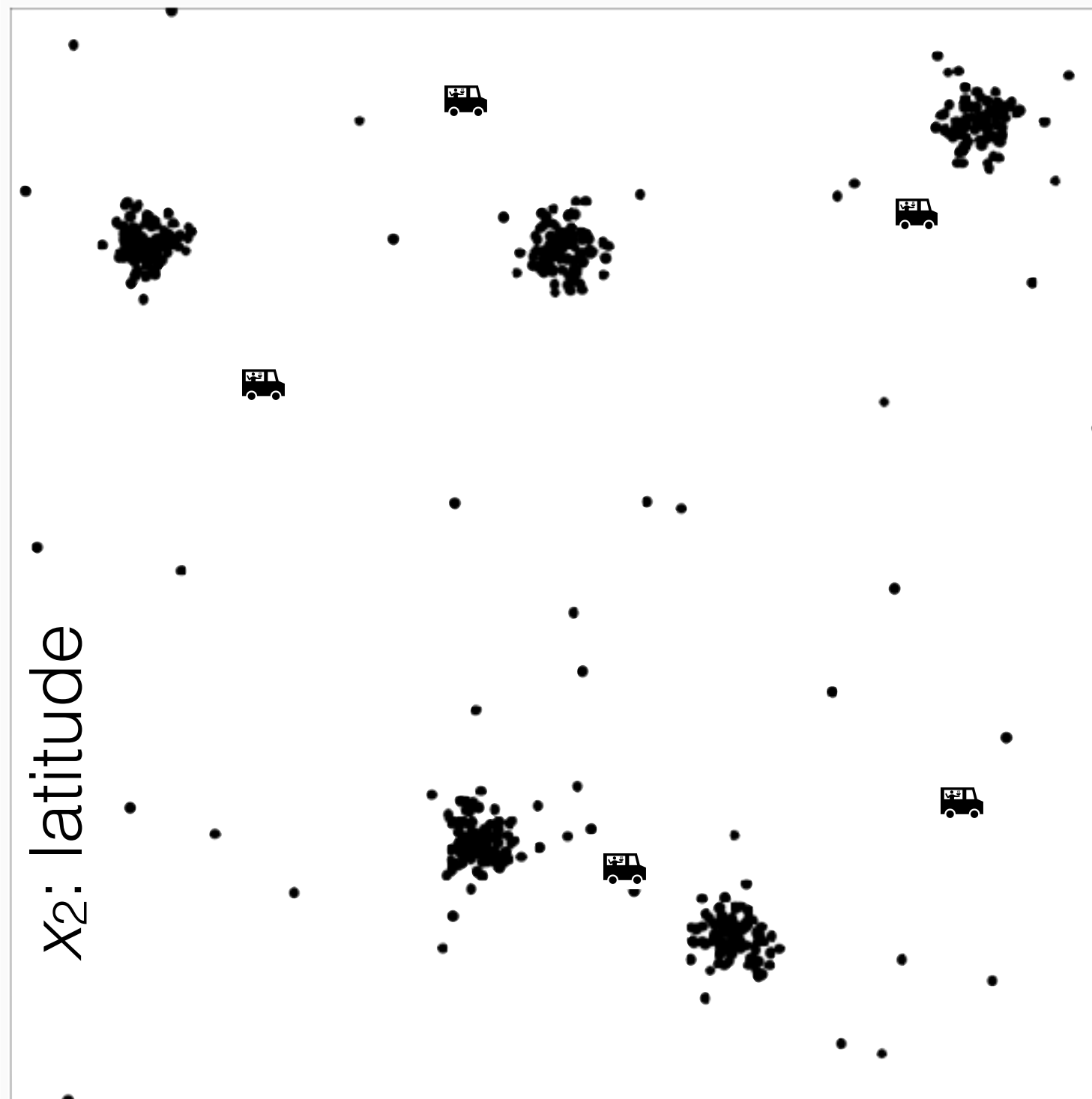$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:
$$\|x^{(i)} - \mu^{(j)}\|_2^2$$
- Loss across all people:

$$\sum_{i=1}^{n} \sum_{j=1}^{k} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

3

# Food distribution placement



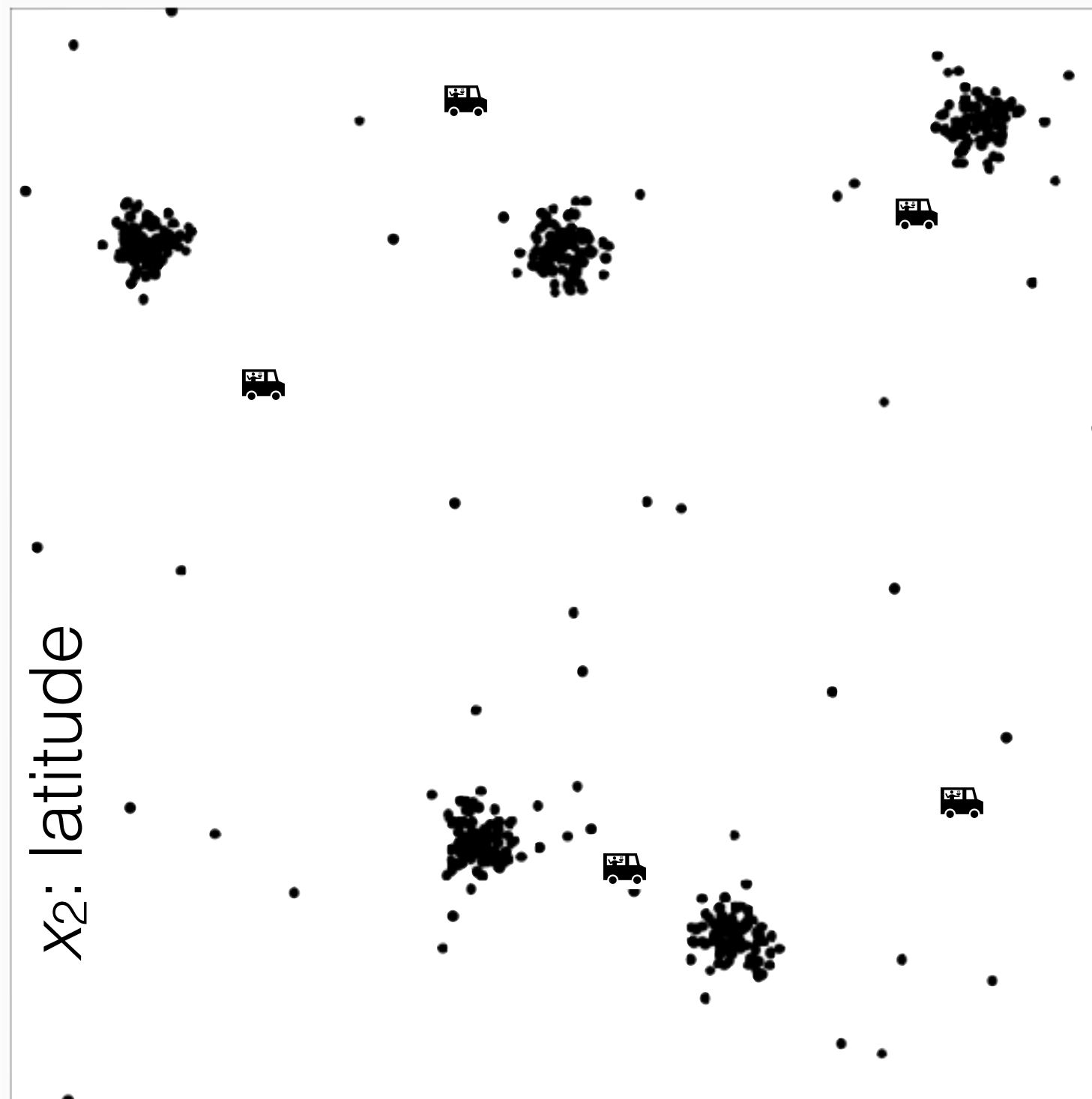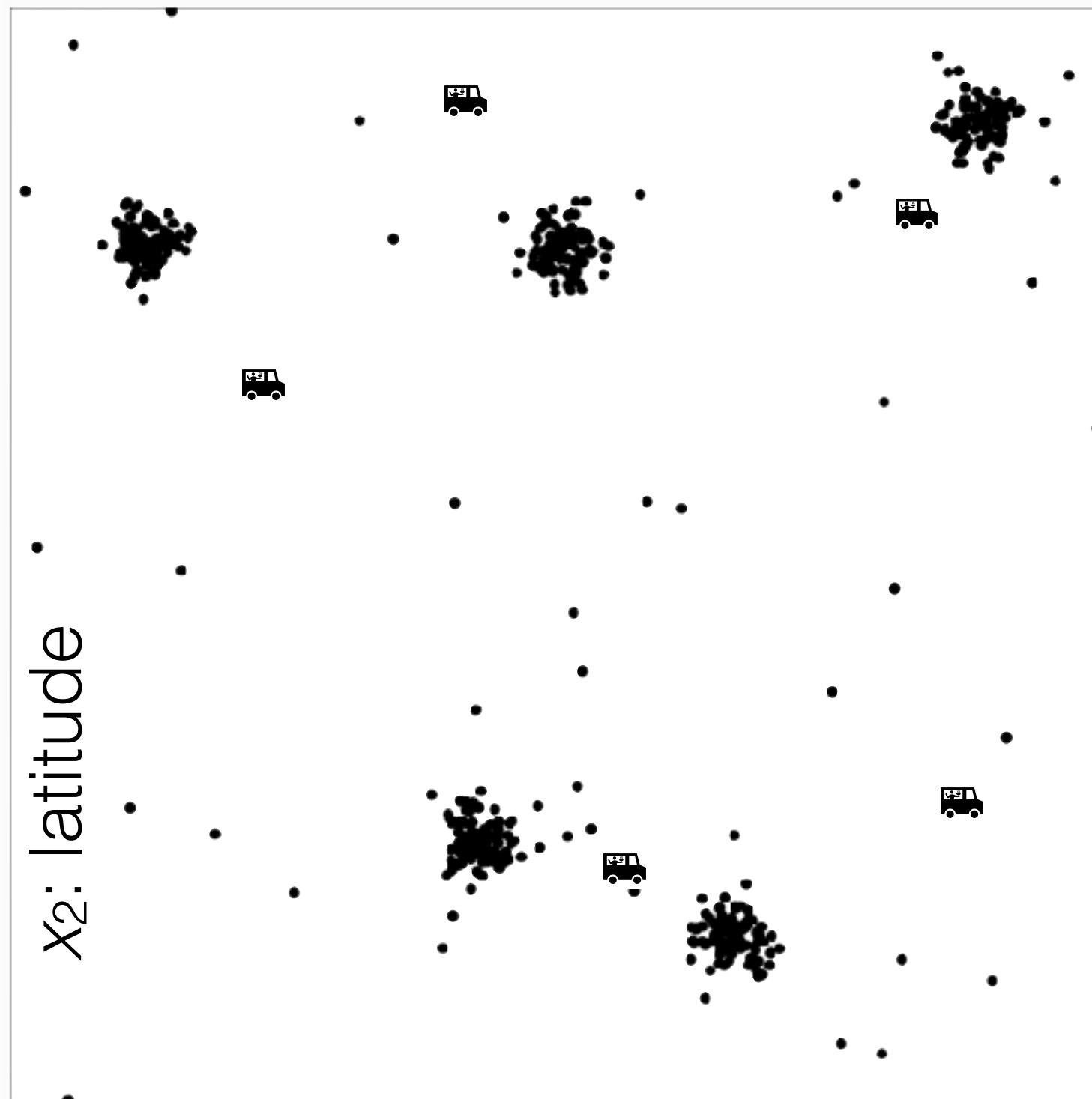$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:
$$\|x^{(i)} - \mu^{(j)}\|_2^2$$
- Loss across all people:

$$\sum_{i=1}^{n} \sum_{j=1}^{k} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

3

# Food distribution placement
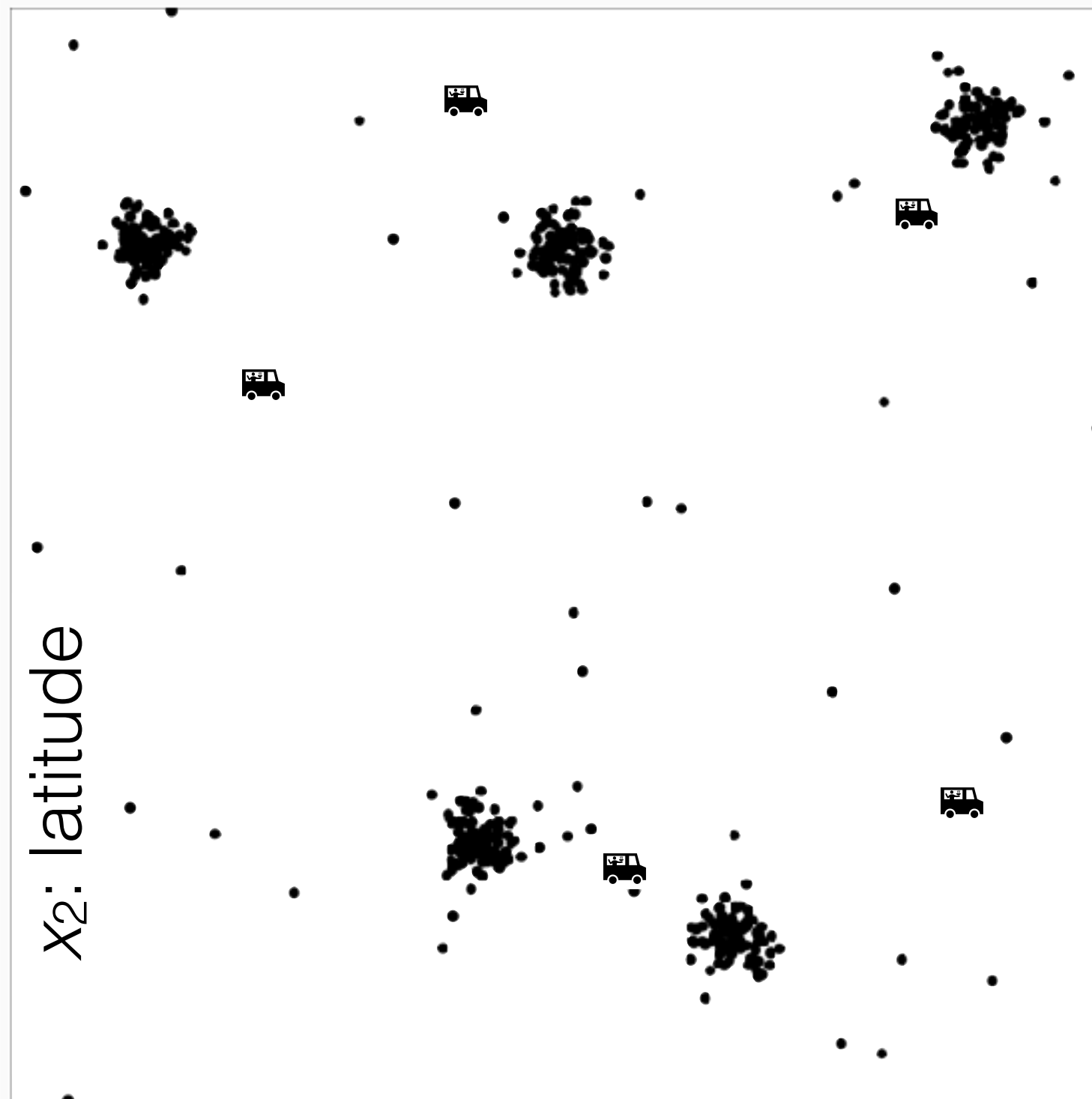


$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:
$$\|x^{(i)} - \mu^{(j)}\|_2^2$$
- Loss across all people:

$$\sum_{i=1}^{n} \sum_{j=1}^{k} \mathbf{1}\{y^{(i)} = j\} \|x^{(i)} - \mu^{(j)}\|_2^2$$

3

# Food distribution placement
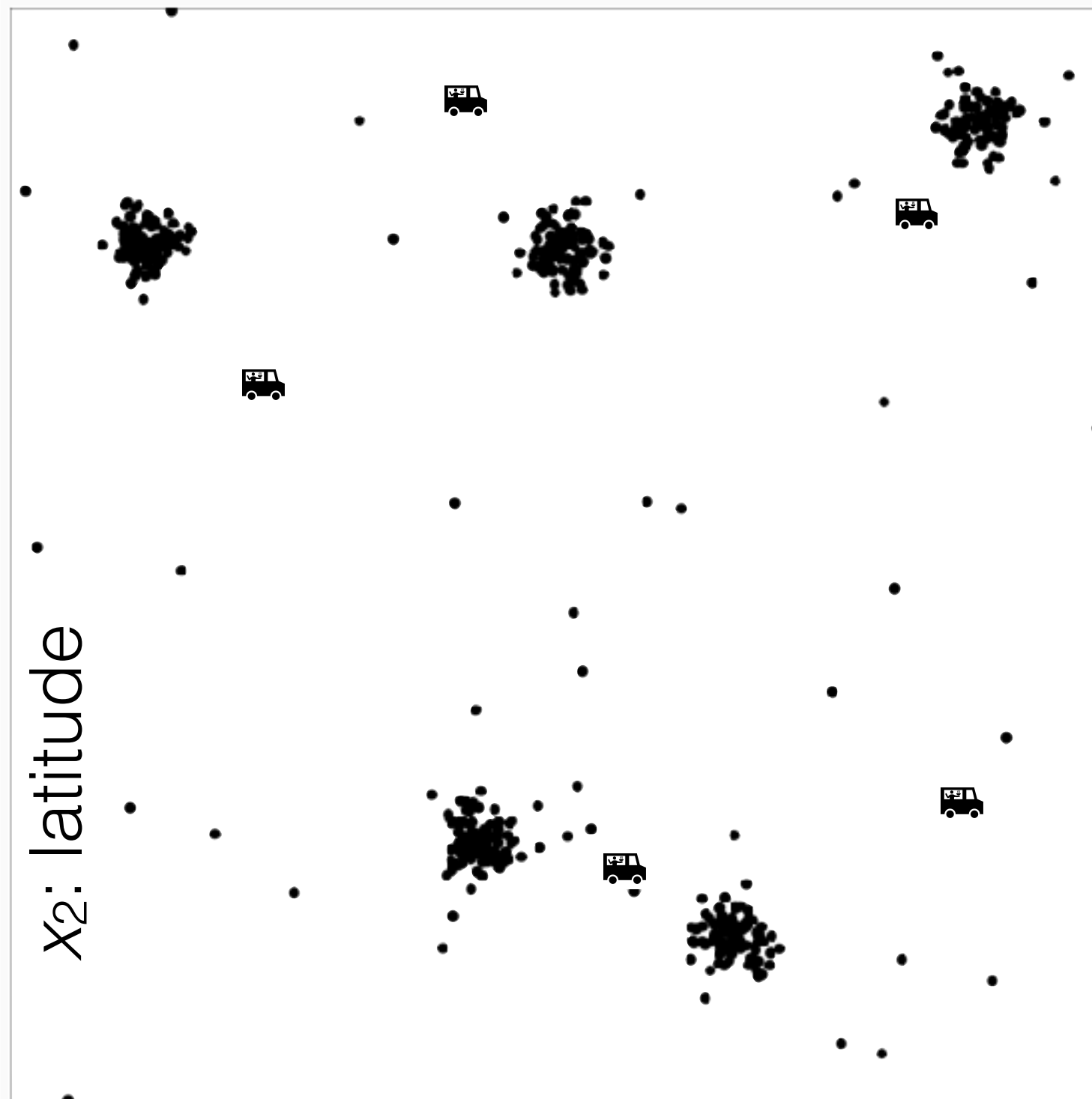
$x_2$: latitude

$x_1$: longitude

- Where should I have my *k* food trucks park?
- Want to minimize the loss of people we serve
- Person *i* location $x^{(i)}$
- Food truck *j* location $\mu^{(j)}$
- Index of truck where person *i* walks: $y^{(i)}$
- Loss if *i* walks to truck *j*:
$$\|x^{(i)} - \mu^{(j)}\|_2^2$$
- Loss across all people:

$$\sum_{i=1}^{n} \sum_{j=1}^{k} \mathbf{1}\{y^{(i)} = j\} \|x^{(i)} - \mu^{(j)}\|_2^2$$

3

# Food distribution placement
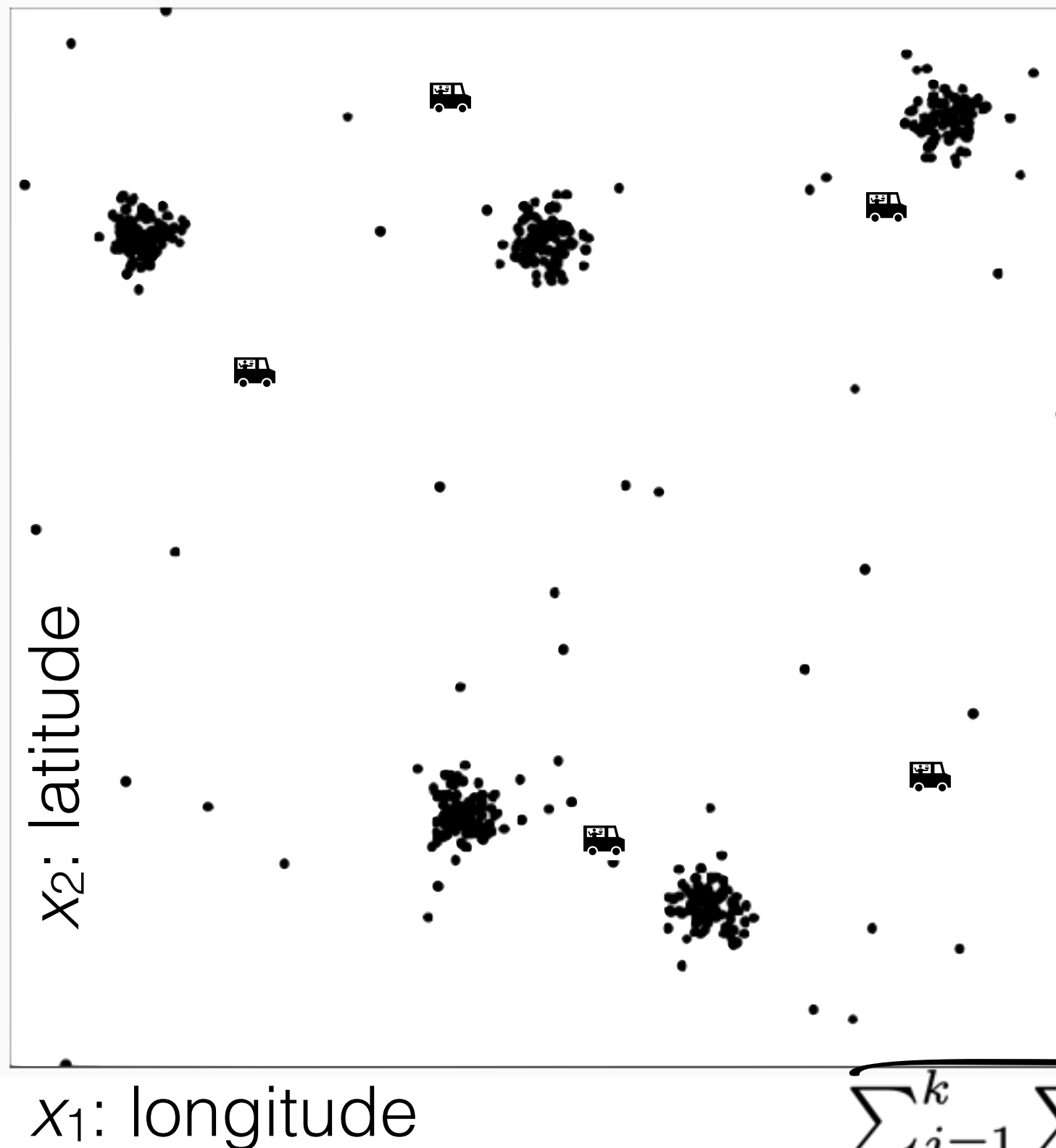


$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:
$$\|x^{(i)} - \mu^{(j)}\|_2^2$$
- Loss across all people:
$$\sum_{i=1}^{n} \sum_{j=1}^{k} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

3

# Food distribution placement
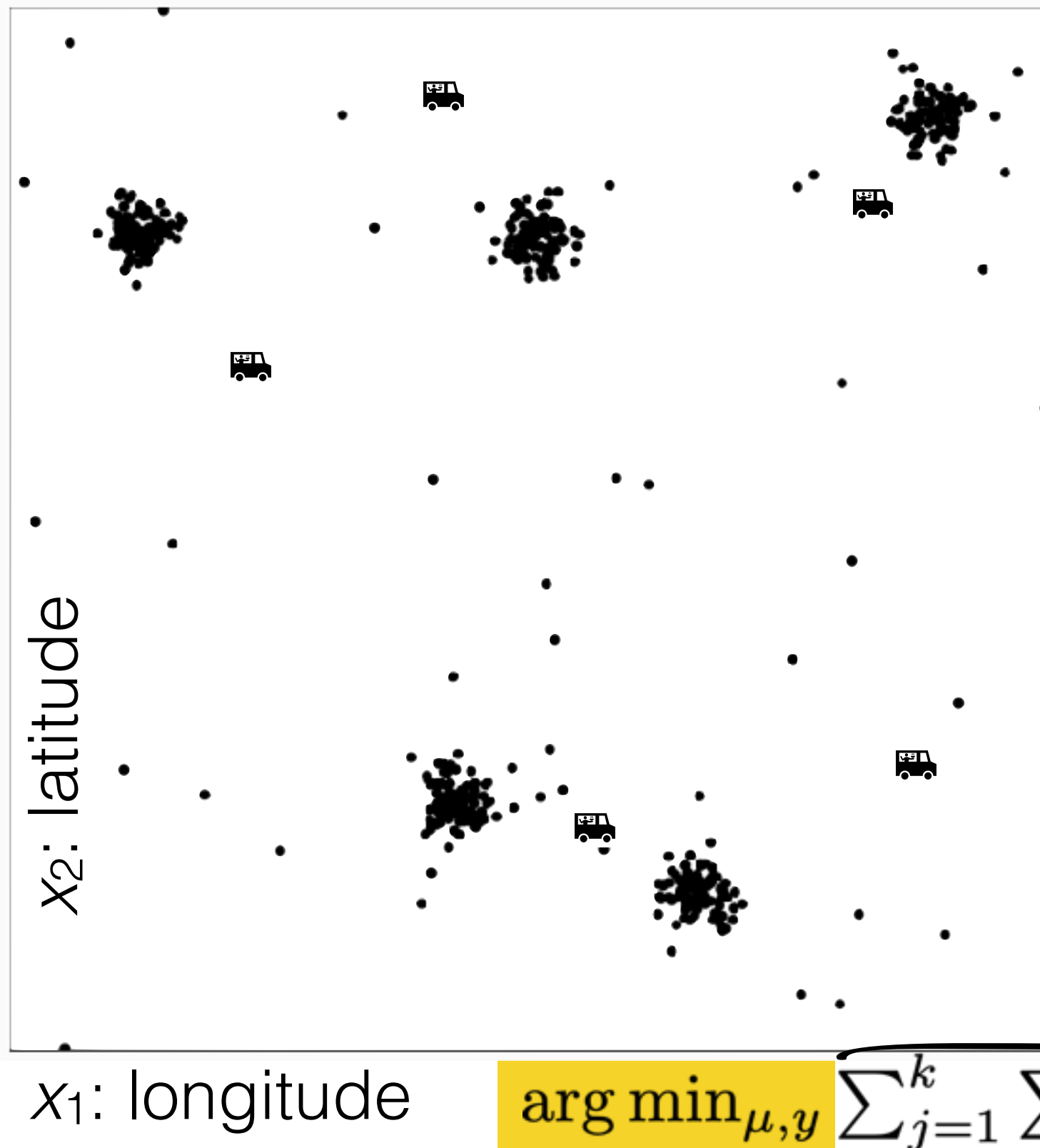
$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:

$$\|x^{(i)} - \mu^{(j)}\|_2^2$$

- Loss across all people:

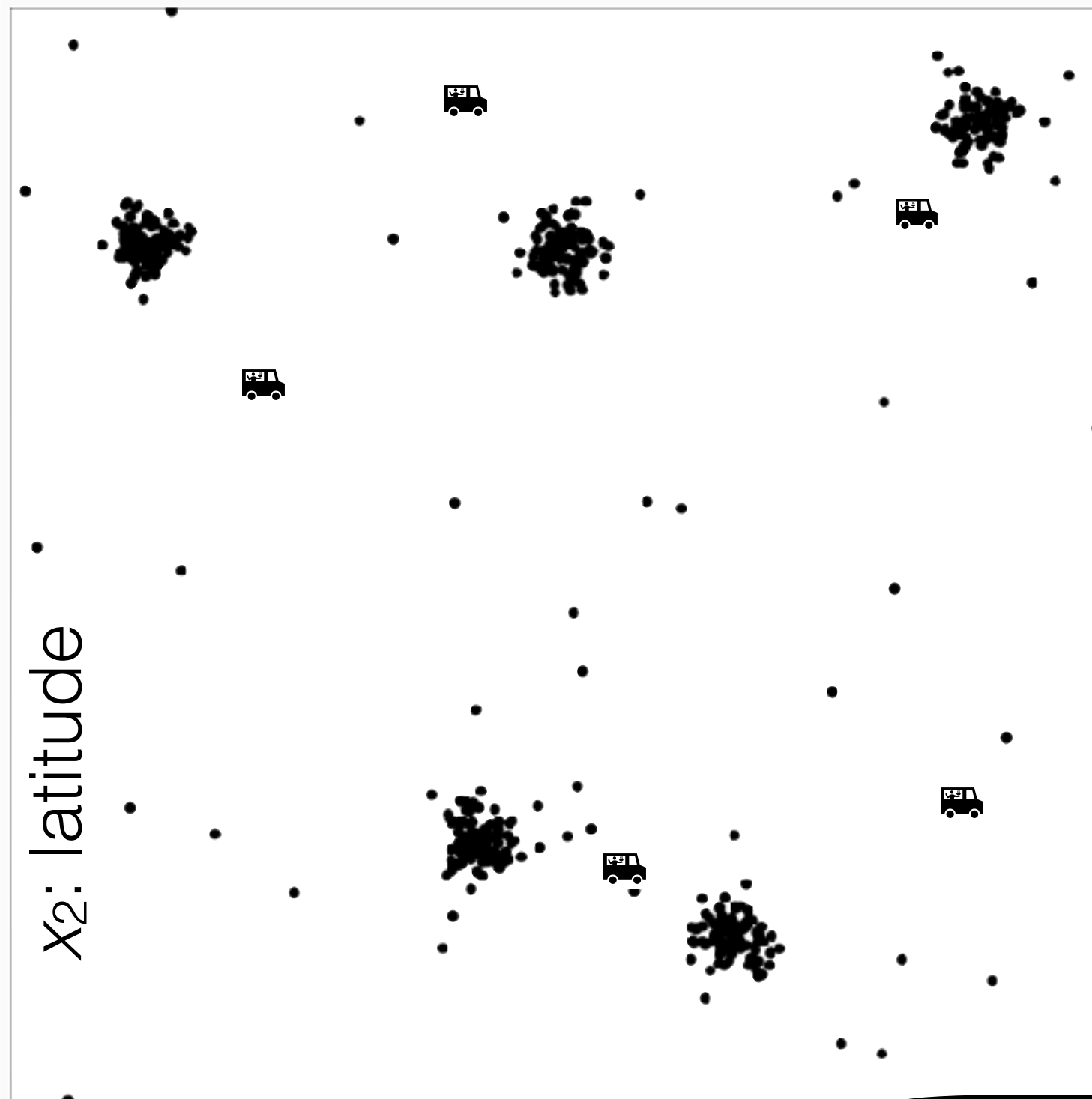$$\sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

3

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]

# Food distribution placement

$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:

$$\|x^{(i)} - \mu^{(j)}\|_2^2$$

- Loss across all people:

$$\sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

3

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]

# Food distribution placement



$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
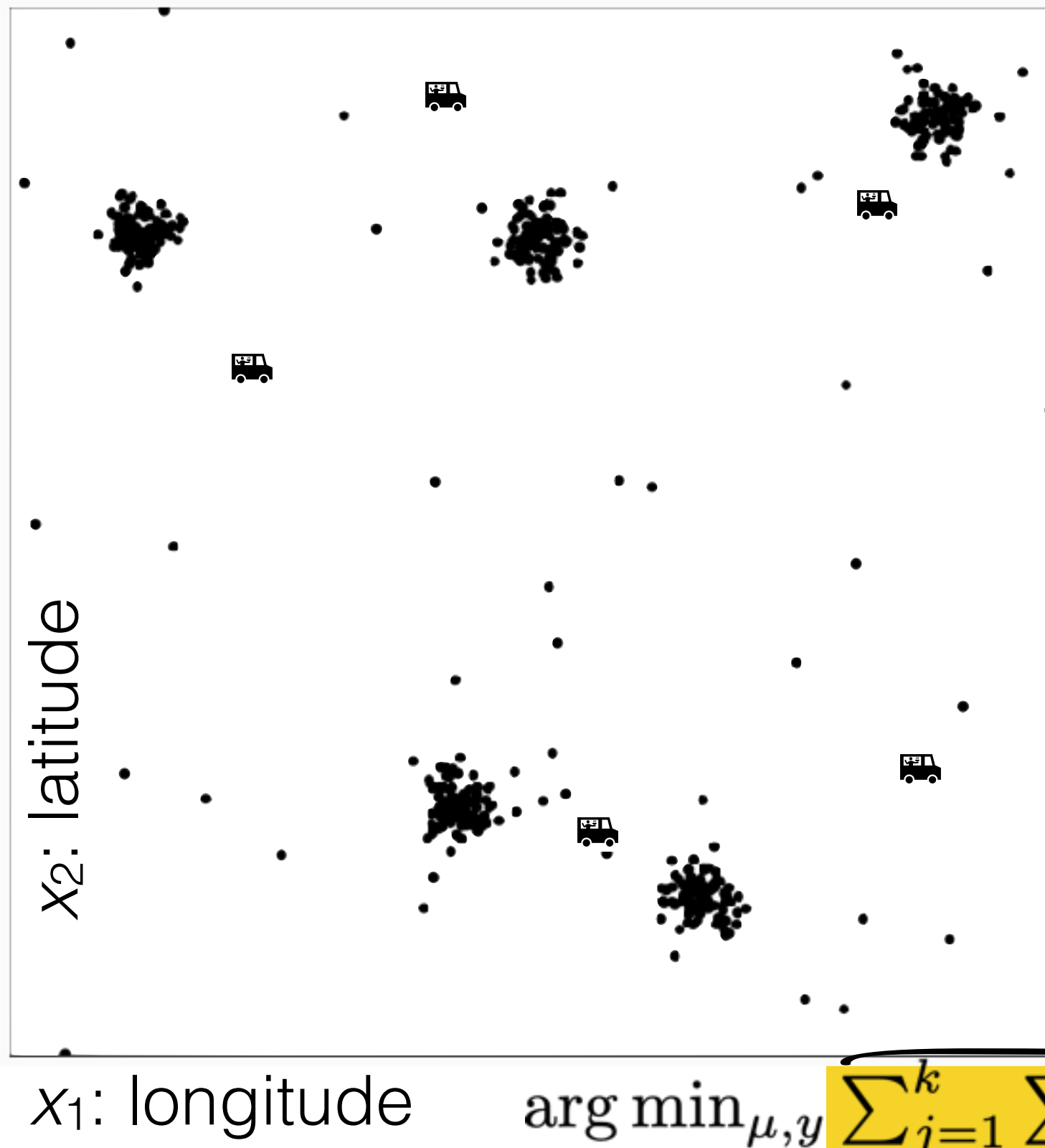- Loss if $i$ walks to truck $j$:
$$\|x^{(i)} - \mu^{(j)}\|_2^2$$
- Loss across all people:

$$\arg\min_{\mu,y} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

3

# Food distribution placement



$x_2$: latitude

$x_1$: longitude

- Where should I have my $k$ food trucks park?
- Want to minimize the loss of people we serve
- Person $i$ location $x^{(i)}$
- Food truck $j$ location $\mu^{(j)}$
- Index of truck where person $i$ walks: $y^{(i)}$
- Loss if $i$ walks to truck $j$:

$$\|x^{(i)} - \mu^{(j)}\|_2^2$$

- Loss across all people:

$$\arg\min_{\mu,y} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

- a.k.a. *k-means objective*

3

# Food distribution placement

$x_2$: latitude

$x_1$: longitude

- Where should I have my *k* food trucks park?
- Want to minimize the loss of people we serve
- Person *i* location $x^{(i)}$
- Food truck *j* location $\mu^{(j)}$
- Index of truck where person *i* walks: $y^{(i)}$
- Loss if *i* walks to truck *j*:
  $$\|x^{(i)} - \mu^{(j)}\|_2^2$$
- Loss across all people:

$$\arg\min_{\mu, y} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\} \|x^{(i)} - \mu^{(j)}\|_2^2$$

- a.k.a. *k-means objective*

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]

# k-means algorithm

# k-means algorithm

`k-means`

# k-means algorithm

$$k\text{-means}(k, \tau)$$

# k-means algorithm

$$\texttt{k-means}(\texttt{k}, \tau)$$

# k-means algorithm

$$\texttt{k-means(k,}\ \tau\texttt{)}$$

# k-means algorithm

$x_2$: latitude

$x_1$: longitude

4

# k-means algorithm

$x_2$: latitude

$x_1$: longitude

Data here is $n$ feature vectors; no labels

4

# k-means algorithm



`k-means(k,`$\tau$`)`

$x_2$: latitude

$x_1$: longitude

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

`k-means(k, `$\tau$`)`

Init $\{\mu^{(j)}\}_{j=1}^{k}$

4

# k-means algorithm



$\texttt{k-means(k,}\tau\texttt{)}$
  $\texttt{Init } \{\mu^{(j)}\}_{j=1}^{k}$

$x_2$: latitude

$x_1$: longitude

# k-means algorithm



```
k-means(k, τ)
  Init {μ^(j)}^k_{j=1}
  for t = 1 to τ
```

$x_2$: latitude

$x_1$: longitude

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k, τ)
  Init {μ^(j)}_{j=1}^k
  for t = 1 to τ
```

**for** i = 1 to n

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```
$\text{Init } \{\mu^{(j)}\}_{j=1}^k$

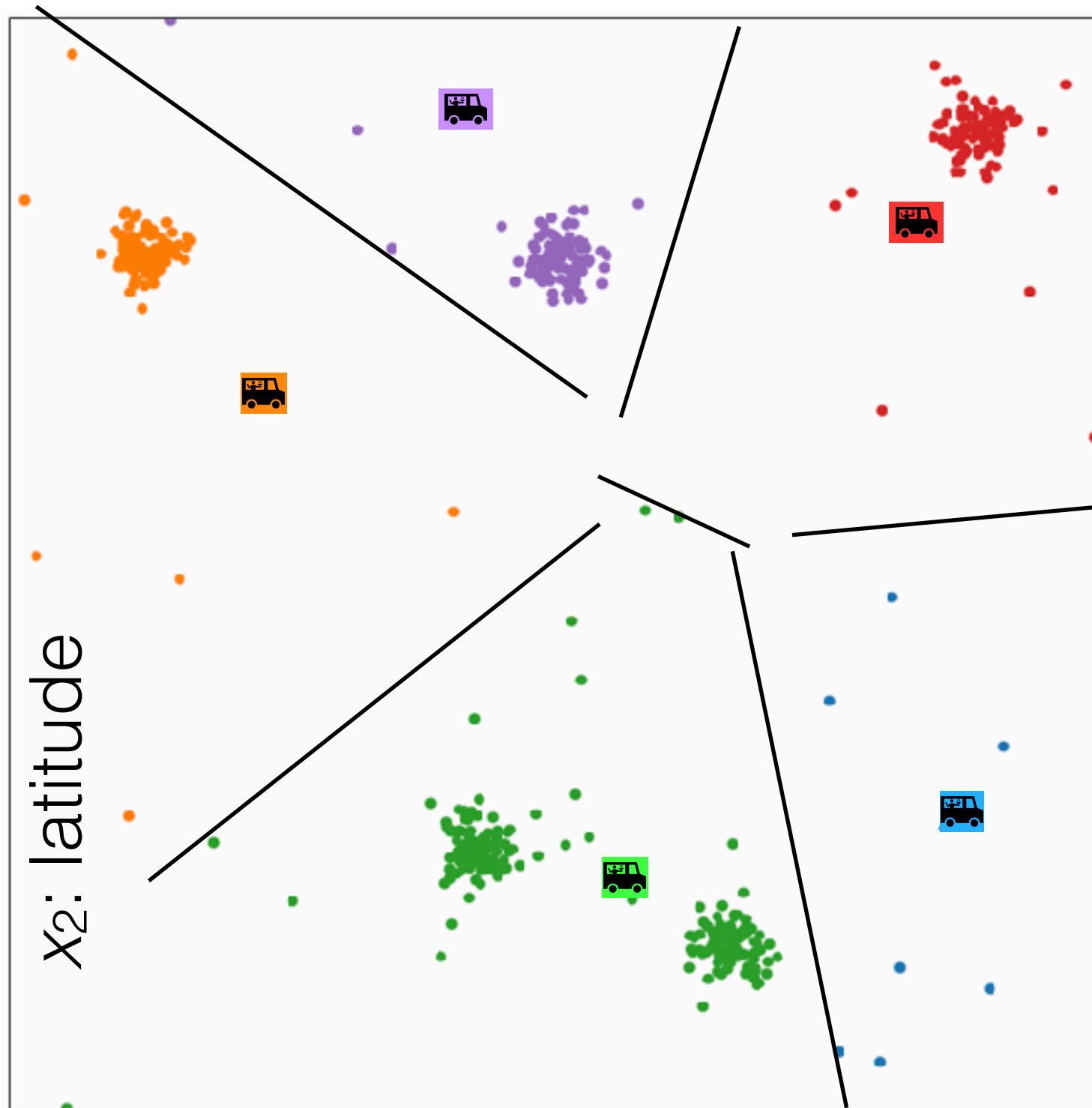**for** t = 1 to $\tau$

    **for** i = 1 to n

$$y^{(i)} = \arg\min_j \|x^{(i)} - \mu^{(j)}\|_2^2$$

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k, τ)
  Init {μ^(j)}_{j=1}^{k}
  for t = 1 to τ

    for i = 1 to n
      y^(i) =
        arg min ||x^(i) − μ^(j)||_2^2
           j
```

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```
$\text{Init } \{\mu^{(j)}\}_{j=1}^{k}$
**for** `t = 1 to` $\tau$

    **for** `i = 1 to n`
$$y^{(i)} = \arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$$

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k, τ)
   Init {μ^(j)}_{j=1}^{k}
   for t = 1 to τ
```

$$\textbf{for } \texttt{i = 1 to n}$$

$$y^{(i)} = \arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$$

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```
Init $\{\mu^{(j)}\}_{j=1}^k$
**for** t = 1 to τ

  **for** i = 1 to n
    $y^{(i)} =$
      $\arg\min_j \|x^{(i)} - \mu^{(j)}\|_2^2$
  **for** j = 1 to k

# k-means algorithm



```
k-means(k,τ)
  Init {μ^(j)}_{j=1}^{k}
  for t = 1 to τ

    for i = 1 to n
      y^(i) =
        arg min_j ‖x^(i) − μ^(j)‖_2^2
    for j = 1 to k
```
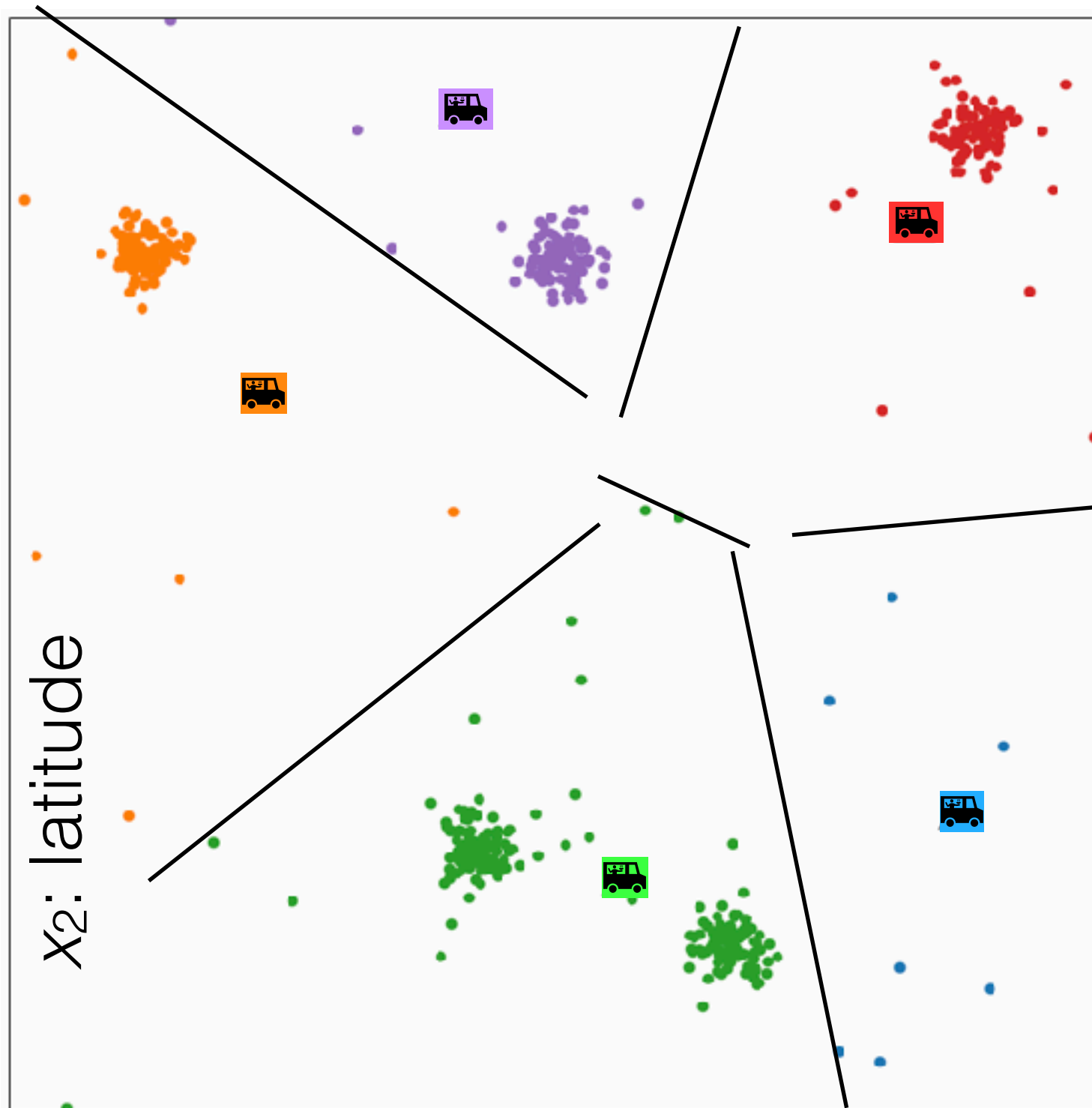
$$\mu^{(j)} = \frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

$x_2$: latitude

$x_1$: longitude

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```
Init $\{\mu^{(j)}\}_{j=1}^{k}$
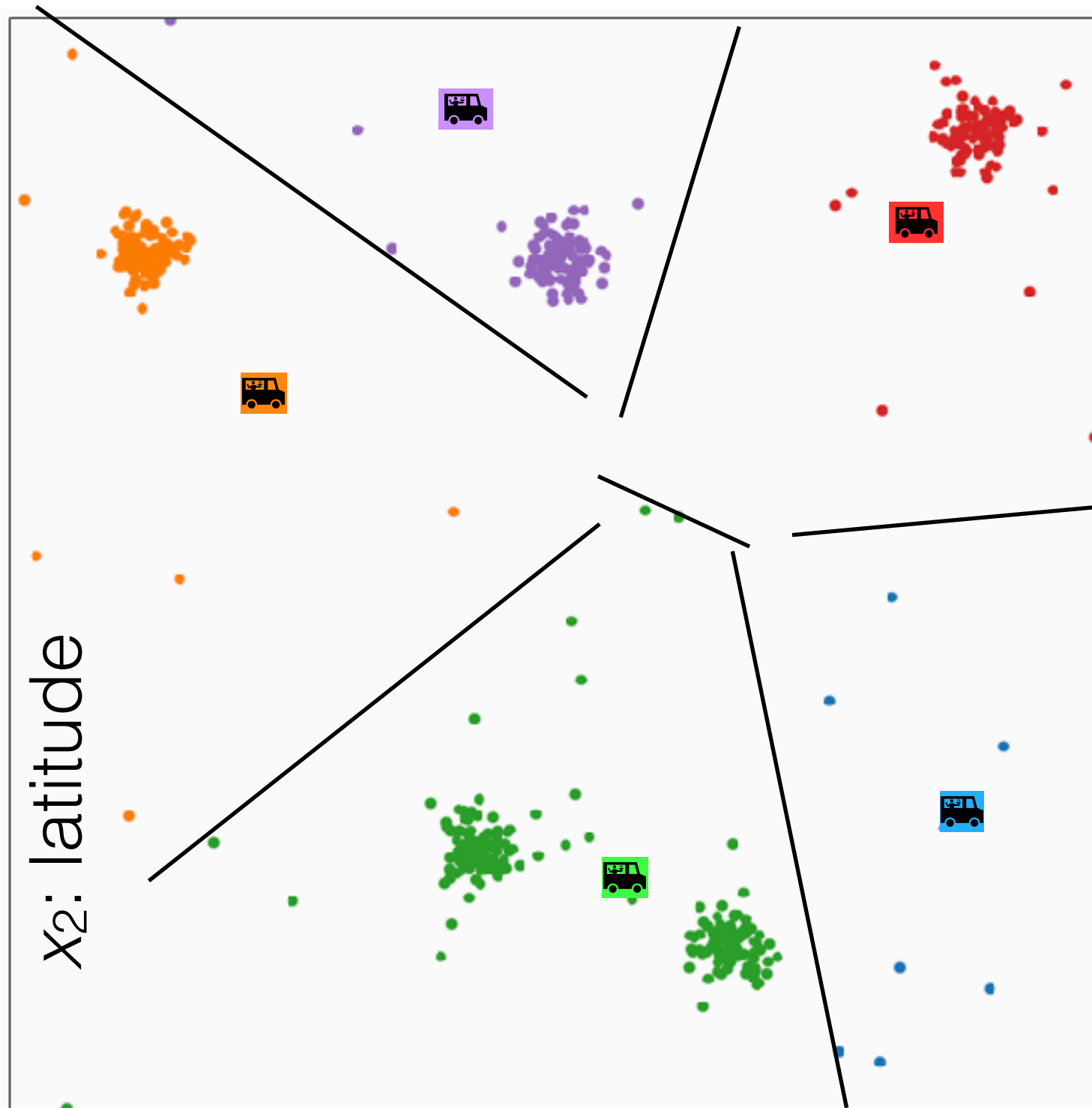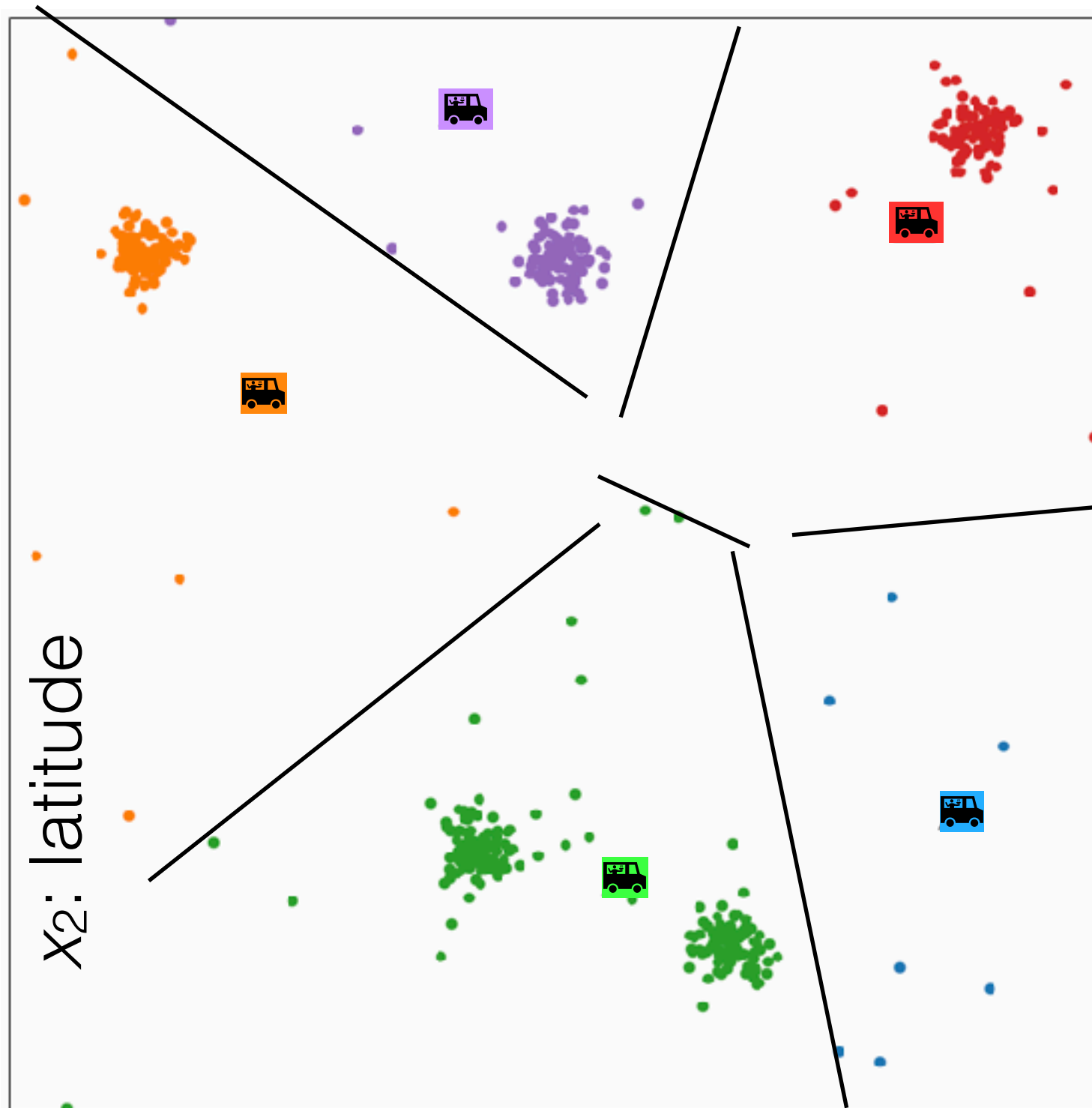**for** t = 1 to τ

    **for** i = 1 to n

$$y^{(i)} = \arg\min_j \|x^{(i)} - \mu^{(j)}\|_2^2$$

    **for** j = 1 to k

$$\mu^{(j)} = \frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k, τ)
```
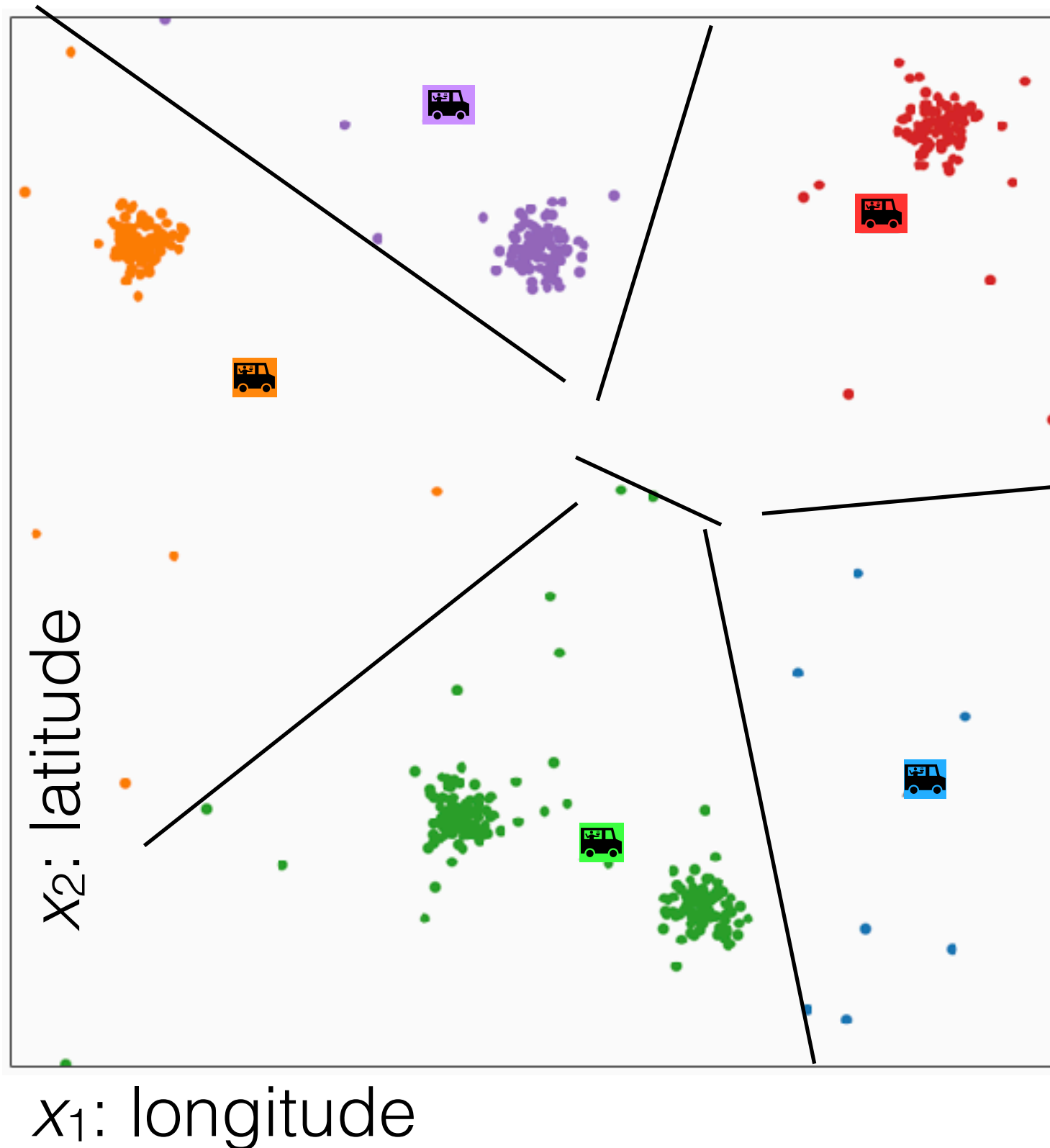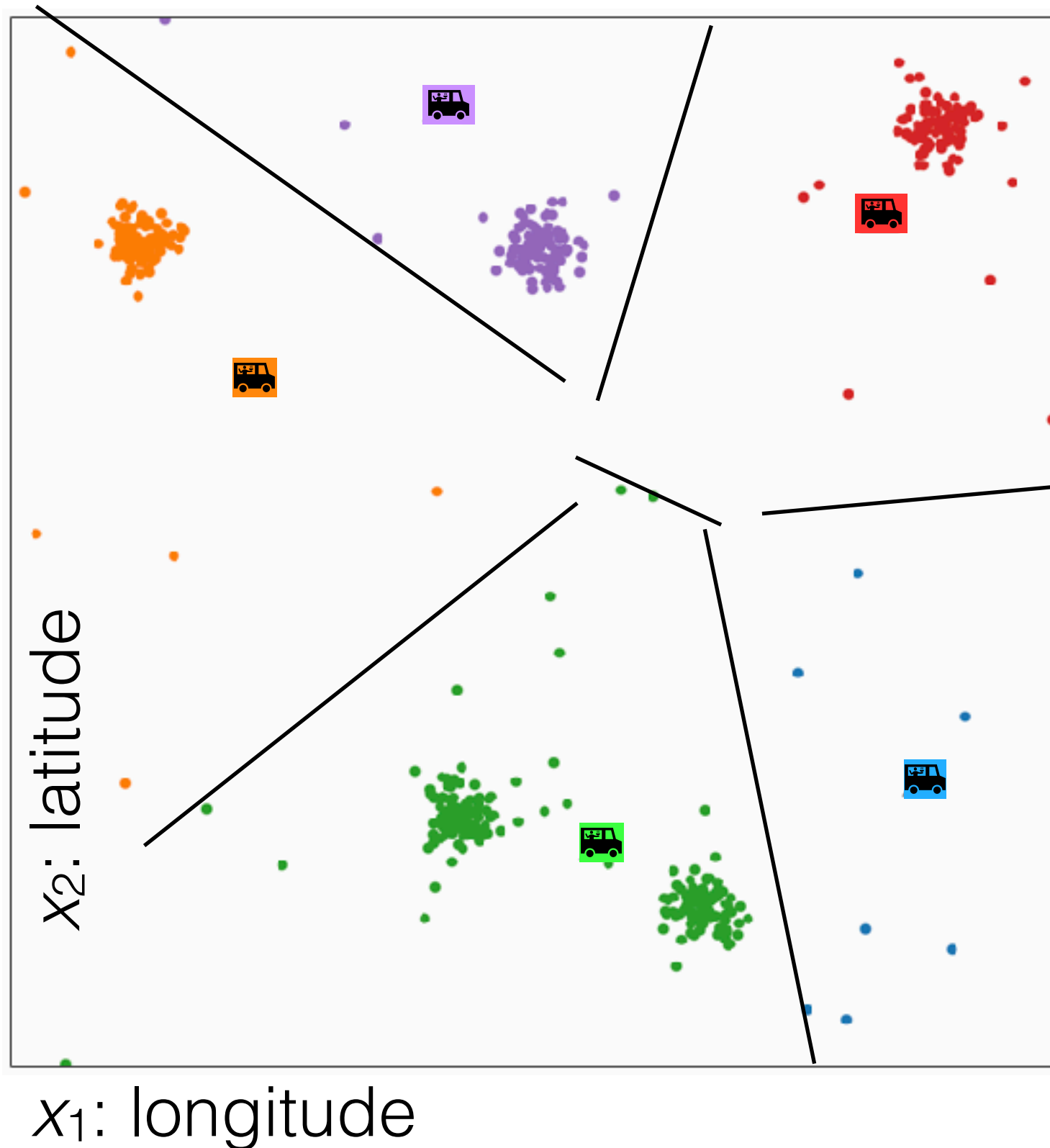Init $\{\mu^{(j)}\}_{j=1}^k$
**for** t = 1 to $\tau$

**for** i = 1 to n
$$y^{(i)} = \arg\min_j \|x^{(i)} - \mu^{(j)}\|_2^2$$
**for** j = 1 to k
$$\mu^{(j)} = \frac{\sum_{i=1}^n \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^n \mathbf{1}\{y^{(i)} = j\}}$$

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```
$$\text{Init } \{\mu^{(j)}\}_{j=1}^{k}$$
**for** t = 1 to τ

    **for** i = 1 to n
$$y^{(i)} =$$
$$\arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$$
    **for** j = 1 to k
$$\mu^{(j)} =$$
$$\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k, τ)
```
$\quad$ Init $\{\mu^{(j)}\}_{j=1}^{k}$

**for** t = 1 to $\tau$

$\quad$ **for** i = 1 to n

$\qquad y^{(i)} =$
$$\arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$$

$\quad$ **for** j = 1 to k

$\qquad \mu^{(j)} =$
$$\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```
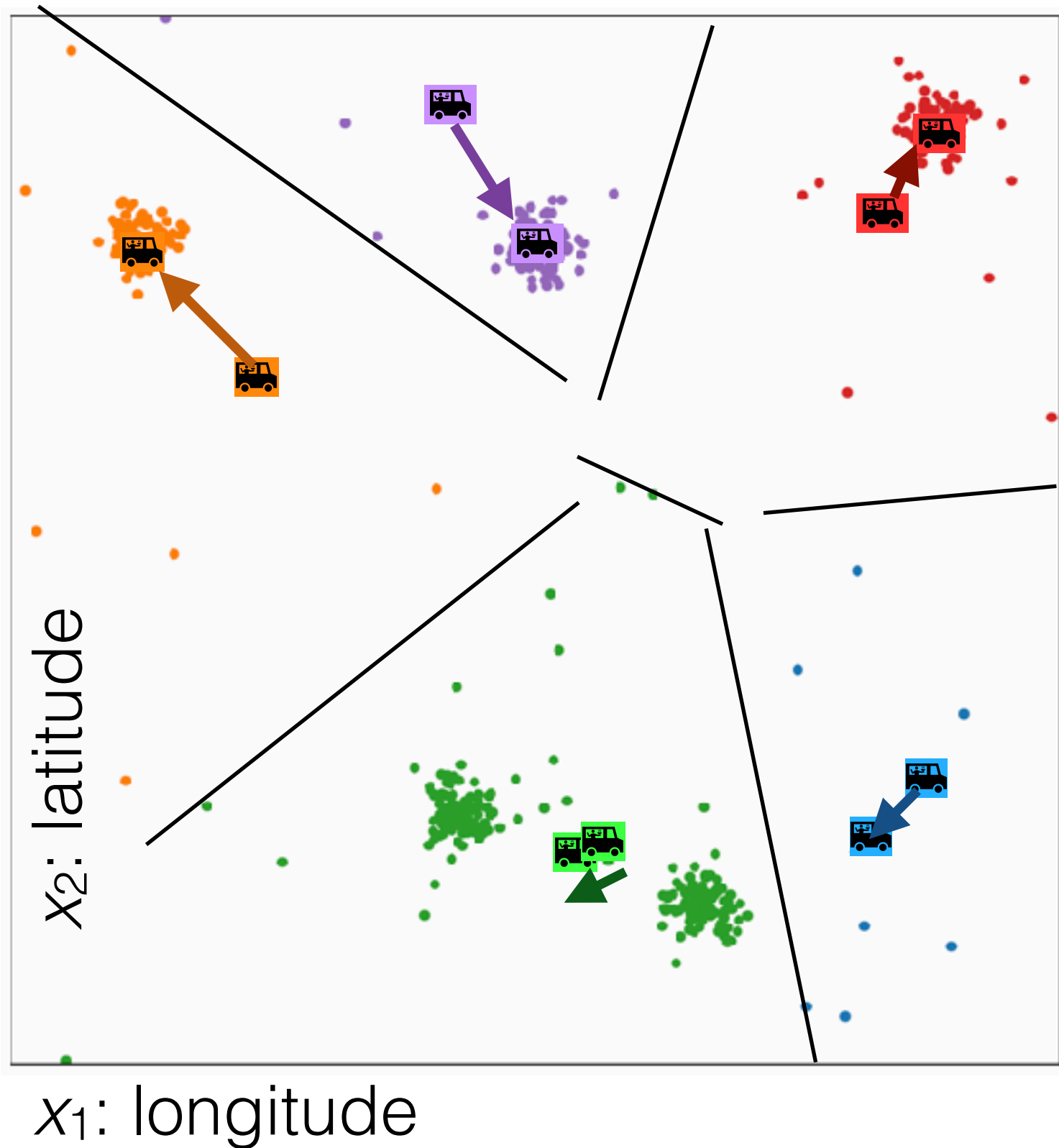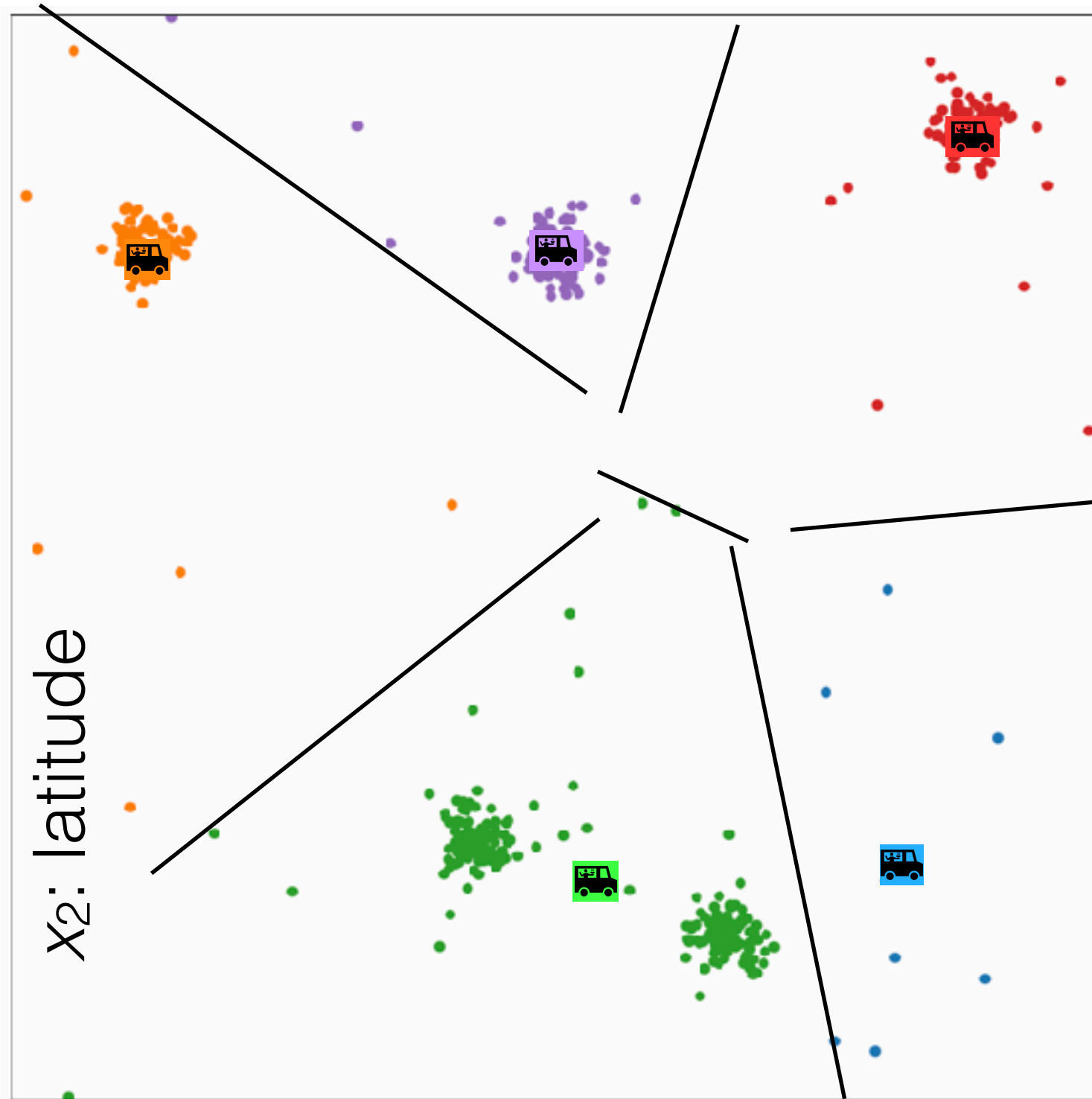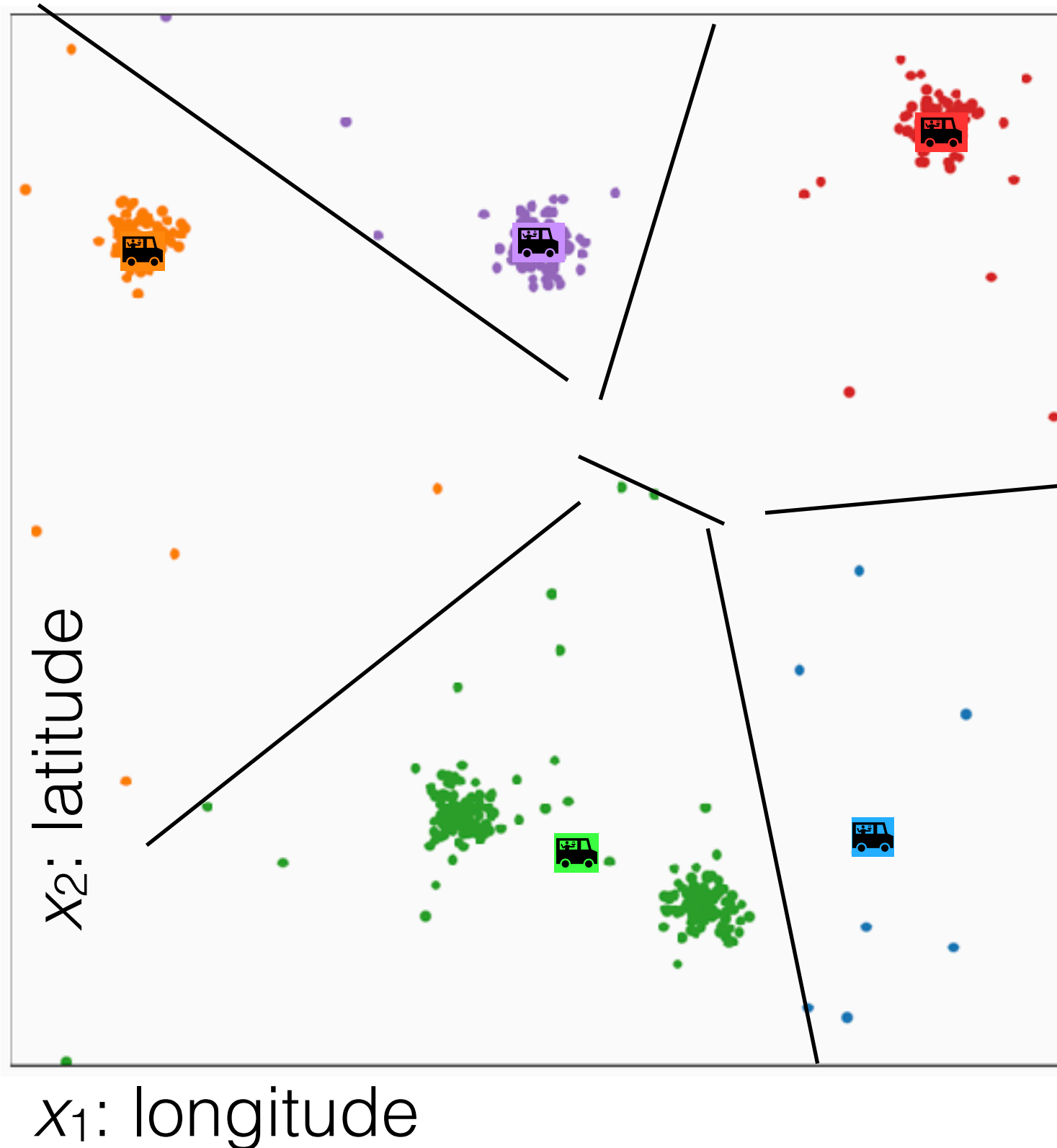$\text{Init } \{\mu^{(j)}\}_{j=1}^{k}$
**for** t = 1 to $\tau$

  **for** i = 1 to n
  $$y^{(i)} = \arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$$
  **for** j = 1 to k
  $$\mu^{(j)} = \frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k, τ)
```
$\quad$ Init $\{\mu^{(j)}\}_{j=1}^{k}$

**for** t = 1 to τ

$\quad$ **for** i = 1 to n

$\qquad y^{(i)} =$
$\qquad \arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$

$\quad$ **for** j = 1 to k

$$\mu^{(j)} = \frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```

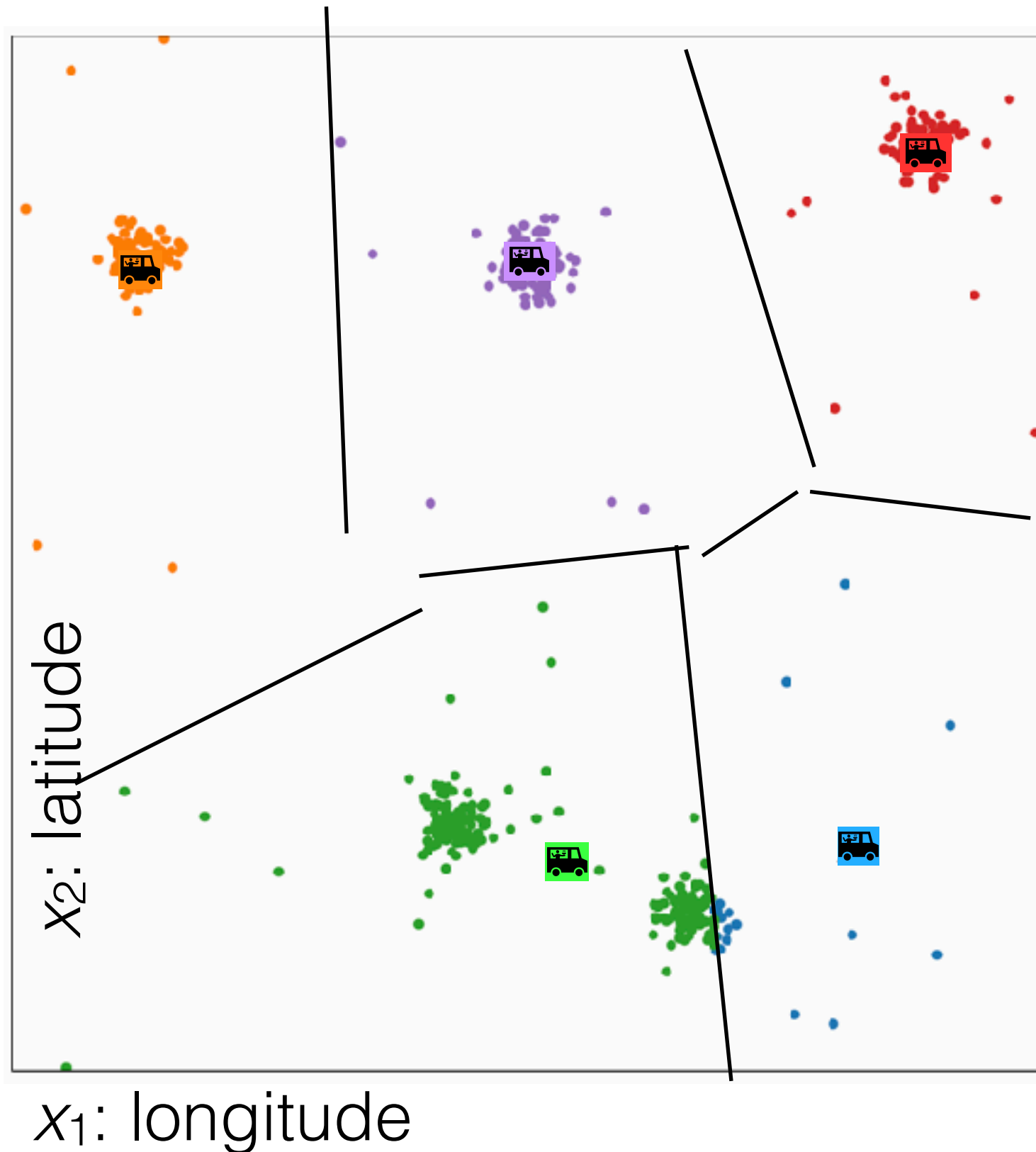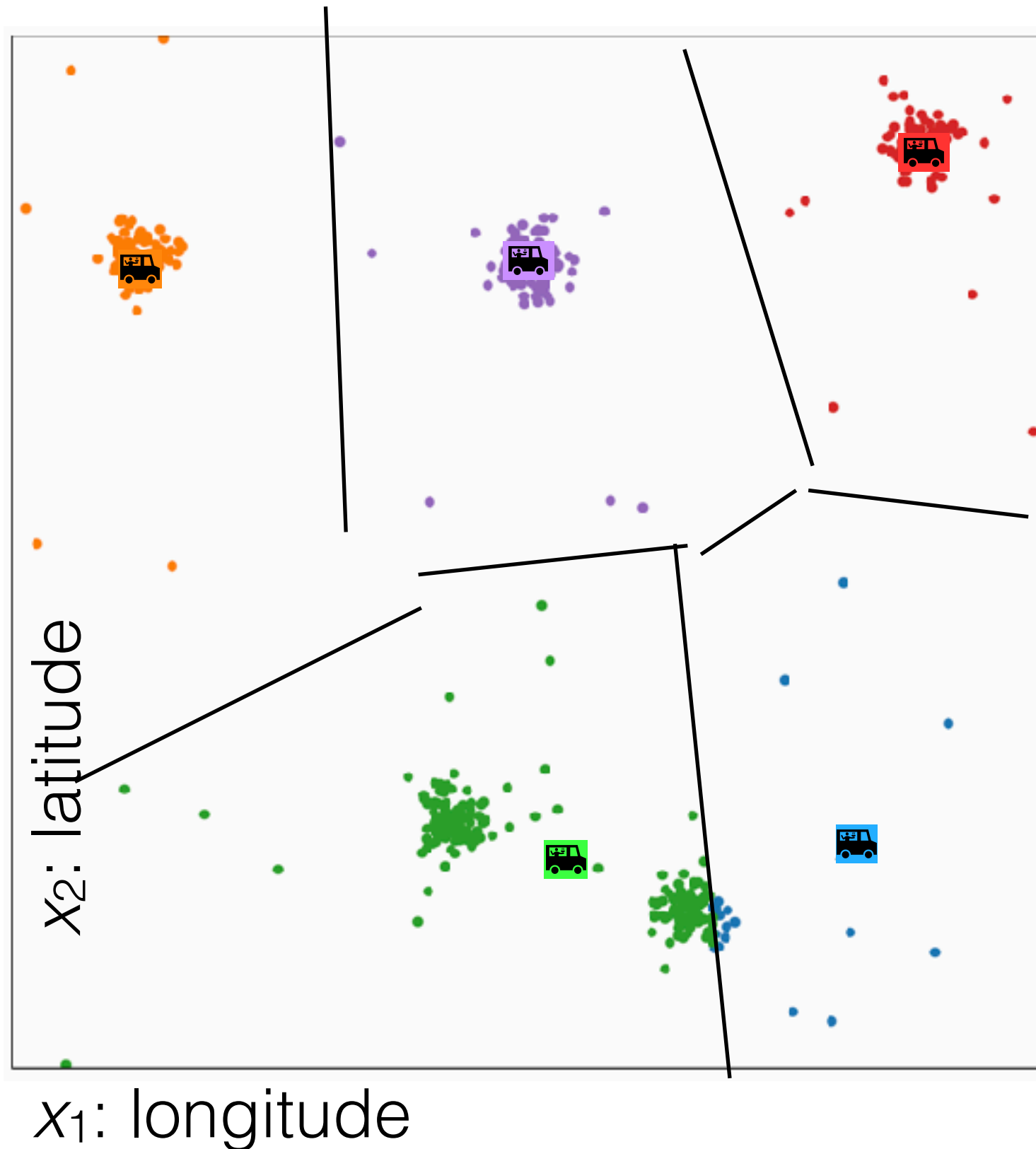$\text{Init } \{\mu^{(j)}\}_{j=1}^{k}$

**for** t = 1 to $\tau$

   **for** i = 1 to n

$$y^{(i)} = \arg\min_j \|x^{(i)} - \mu^{(j)}\|_2^2$$

   **for** j = 1 to k

$$\mu^{(j)} = \frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```
Init $\{\mu^{(j)}\}_{j=1}^{k}$

**for** t = 1 to $\tau$

    **for** i = 1 to n

      $y^{(i)} = $
$$\arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$$

    **for** j = 1 to k

      $\mu^{(j)} = $
$$\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k, τ)
```

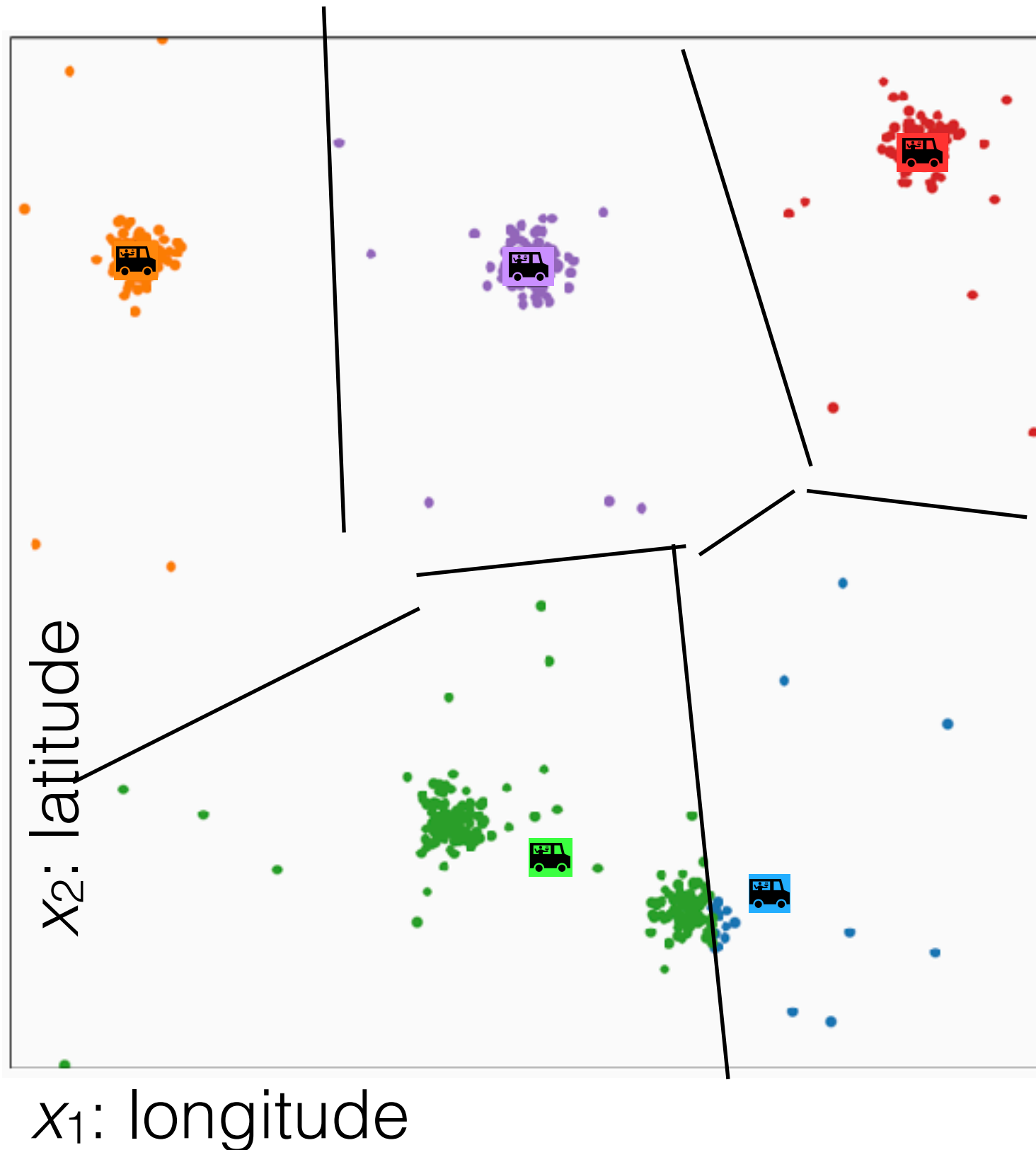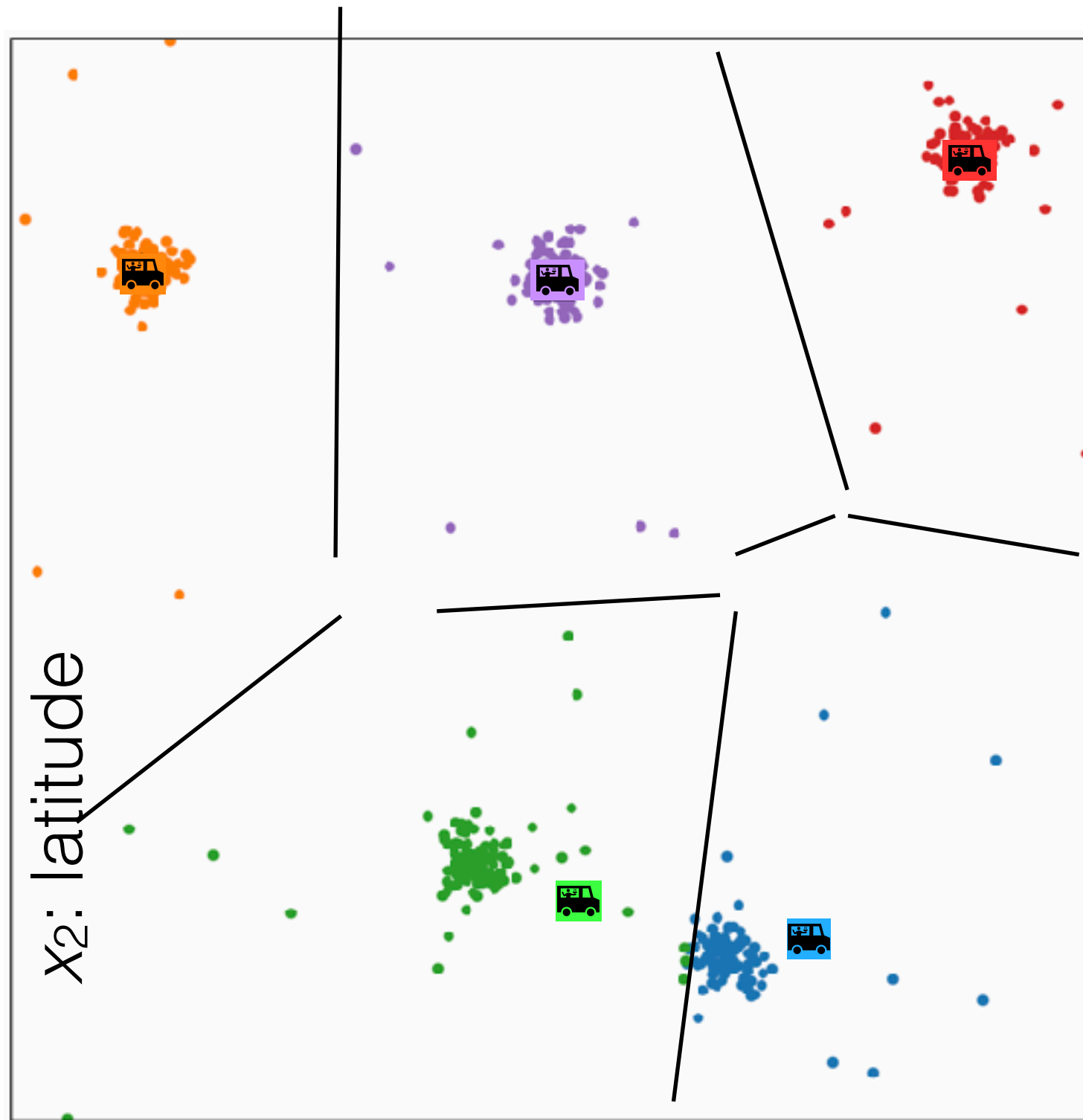Init $\{\mu^{(j)}\}_{j=1}^{k}$

**for** t = 1 to $\tau$

**for** i = 1 to n

$y^{(i)} = \arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$

**for** j = 1 to k

$\mu^{(j)} = \dfrac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k, τ)
```
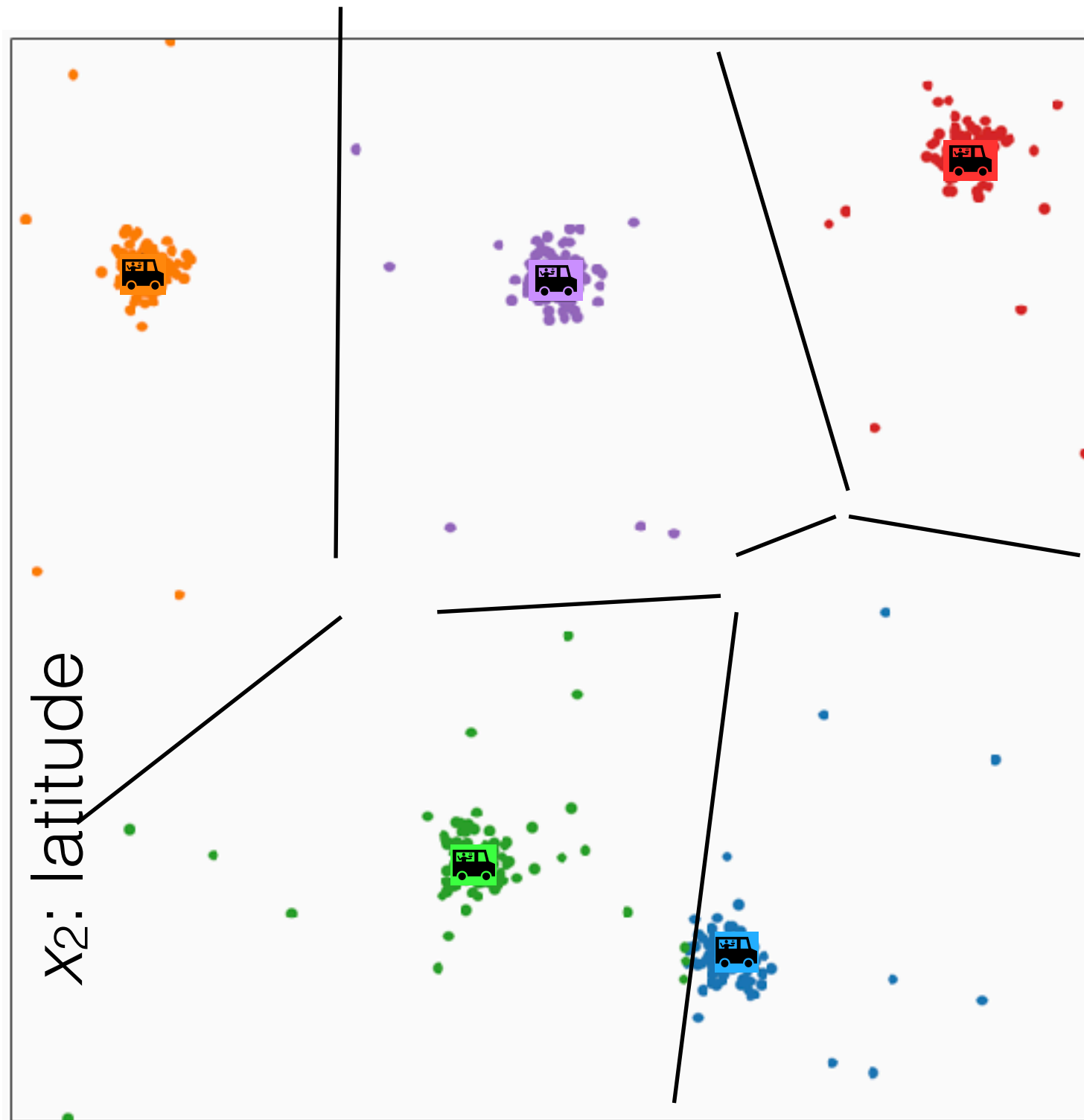Init $\{\mu^{(j)}\}_{j=1}^{k}$

**for** t = 1 to τ

**for** i = 1 to n

$y^{(i)} = \arg\min_j \|x^{(i)} - \mu^{(j)}\|_2^2$

**for** j = 1 to k

$\mu^{(j)} = \dfrac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```
Init $\{\mu^{(j)}\}_{j=1}^{k}$
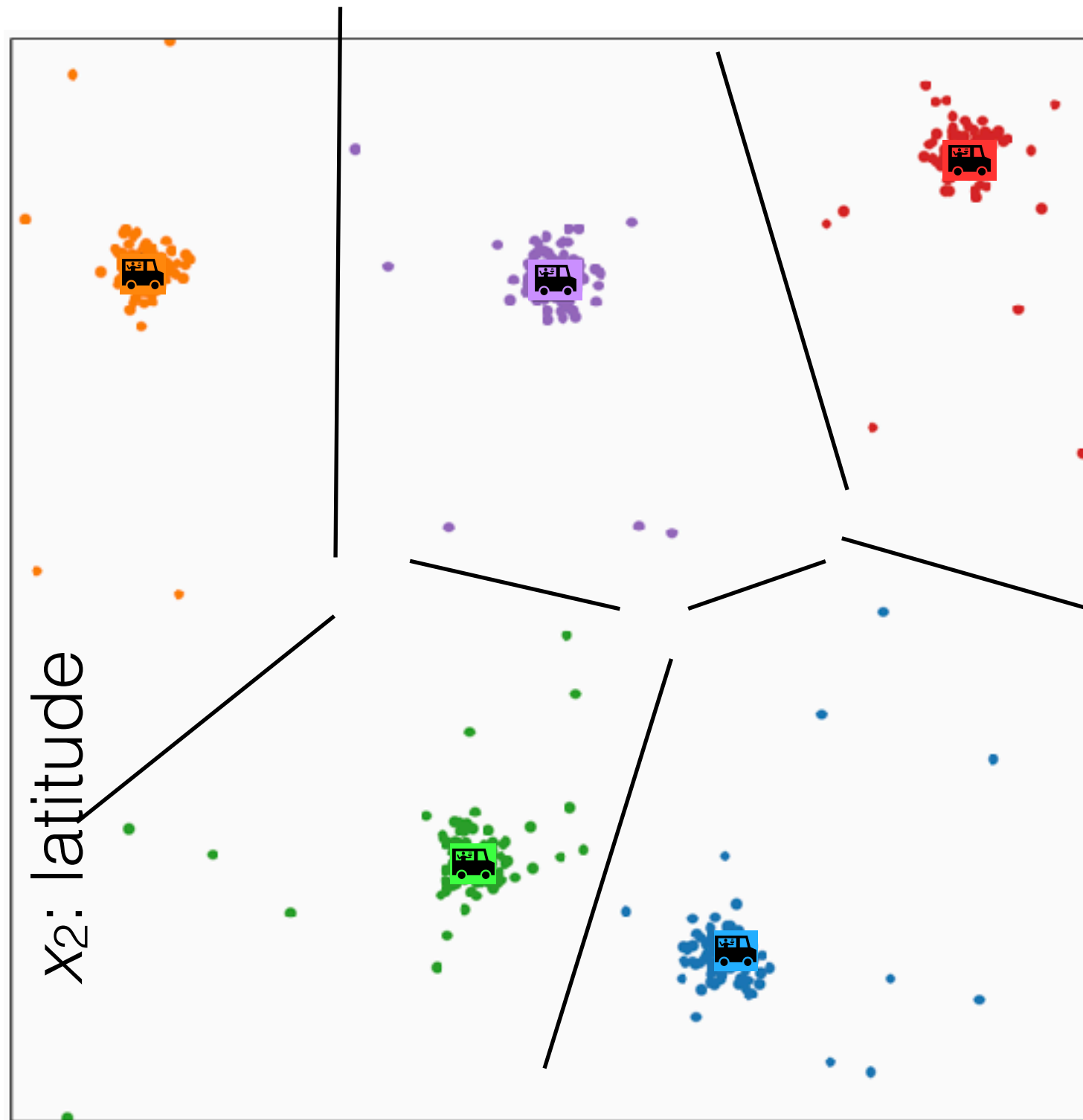**for** t = 1 to $\tau$

   **for** i = 1 to n
   $y^{(i)} =$
   $\arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$

   **for** j = 1 to k
   $\mu^{(j)} =$
   $$\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```
Init $\{\mu^{(j)}\}_{j=1}^{k}$

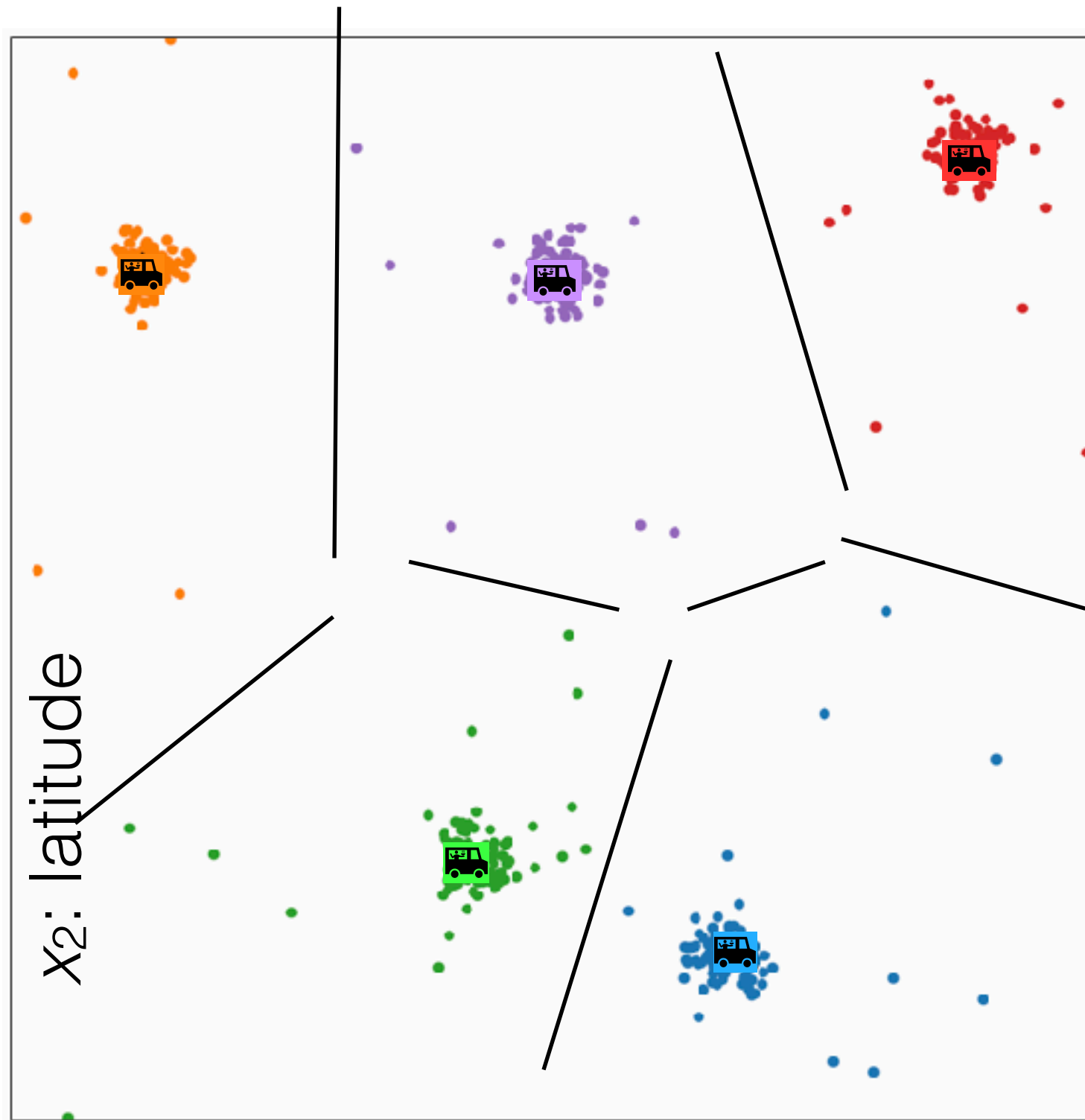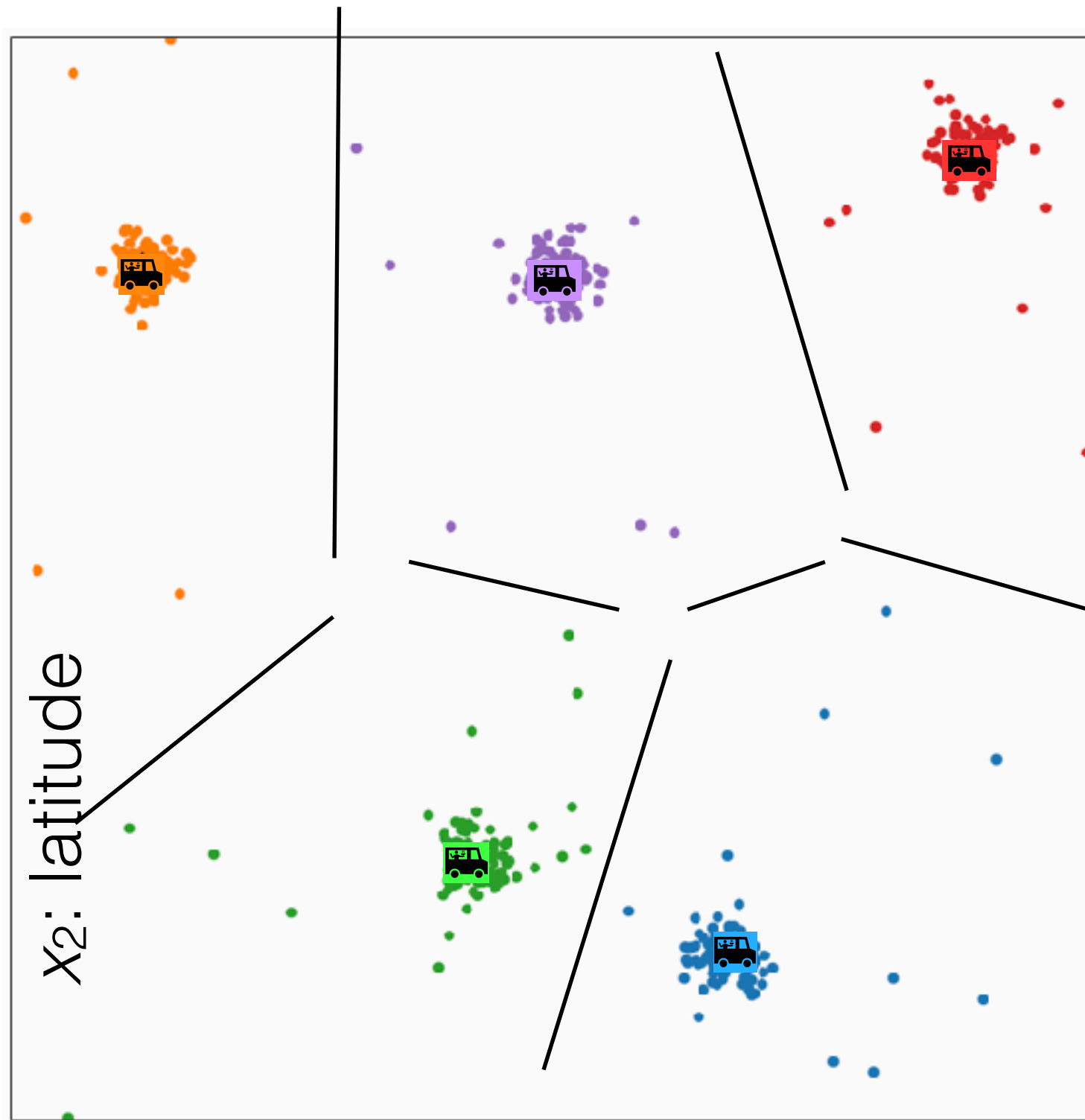**for** t = 1 to $\tau$

    **for** i = 1 to n

    $y^{(i)} =$
$$\arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$$

    **for** j = 1 to k

    $\mu^{(j)} =$
$$\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

# k-means algorithm



```
k-means(k,τ)
```
$\text{Init } \{\mu^{(j)}\}_{j=1}^{k}$
**for** t = 1 to $\tau$

  **for** i = 1 to n
  $y^{(i)} =$
  $\displaystyle\arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$
  **for** j = 1 to k
  $\mu^{(j)} =$
  $$\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

$x_2$: latitude

$x_1$: longitude

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
  Init {μ^(j)}_{j=1}^{k}
  for t = 1 to τ

    for i = 1 to n
      y^(i) =
        arg min_j ‖x^(i) − μ^(j)‖_2^2
    for j = 1 to k
      μ^(j) =
        (∑_{i=1}^{n} 1{y^(i) = j}x^(i)) / (∑_{i=1}^{n} 1{y^(i) = j})
```

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
  Init {μ^(j)}^k_{j=1}
  for t = 1 to τ

    for i = 1 to n
```
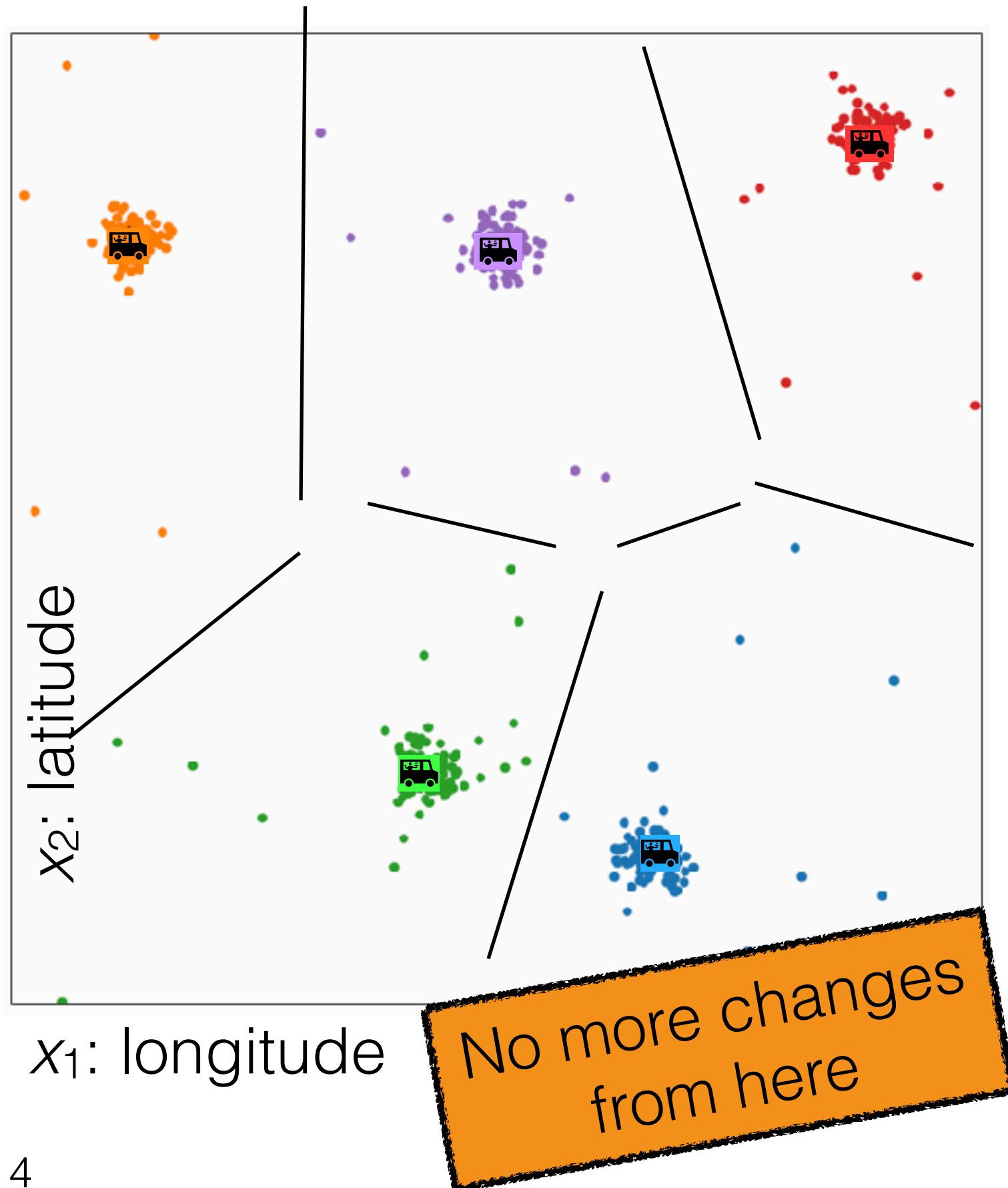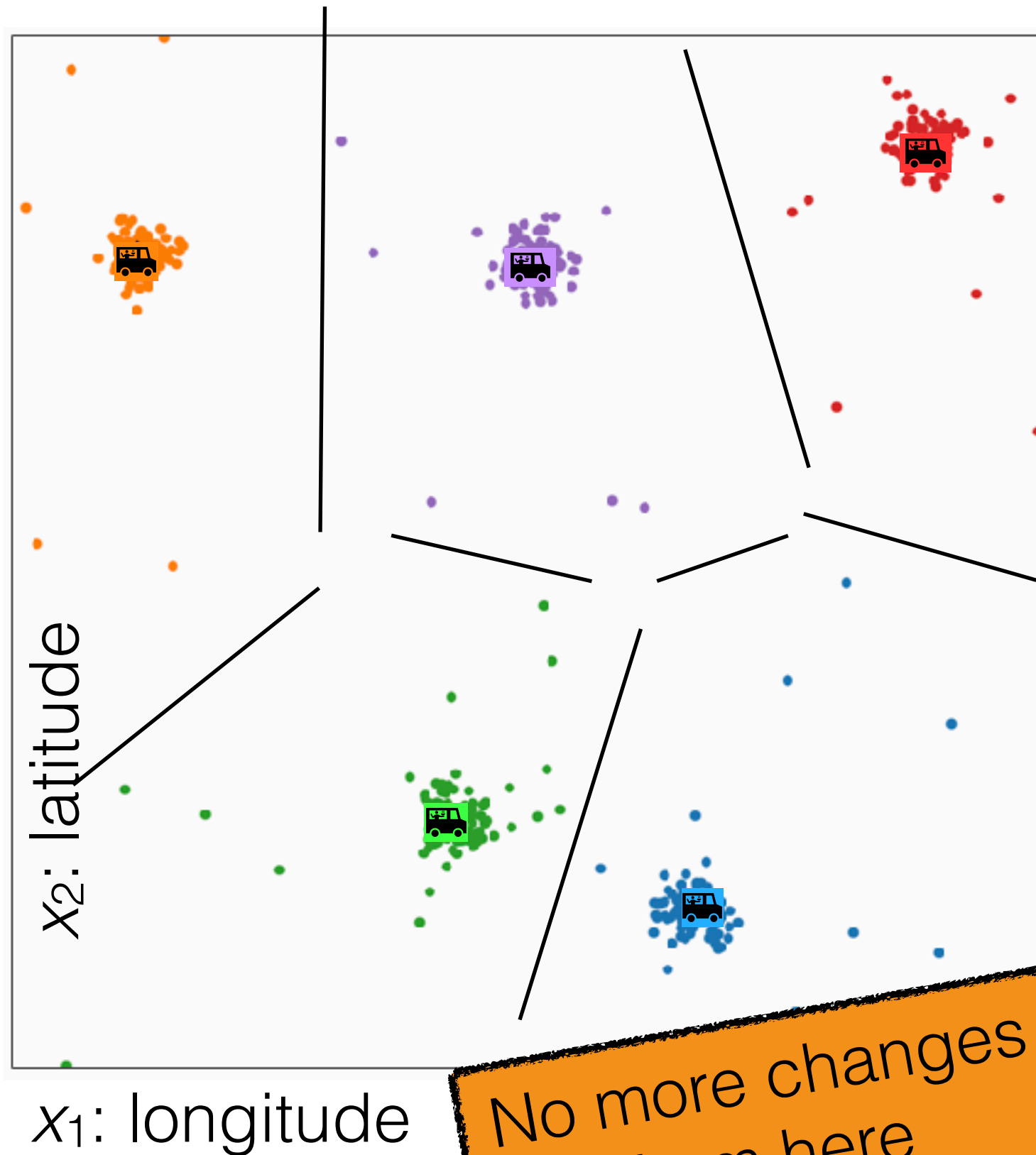$$y^{(i)} = \arg\min_j \|x^{(i)} - \mu^{(j)}\|_2^2$$
```
    for j = 1 to k
```
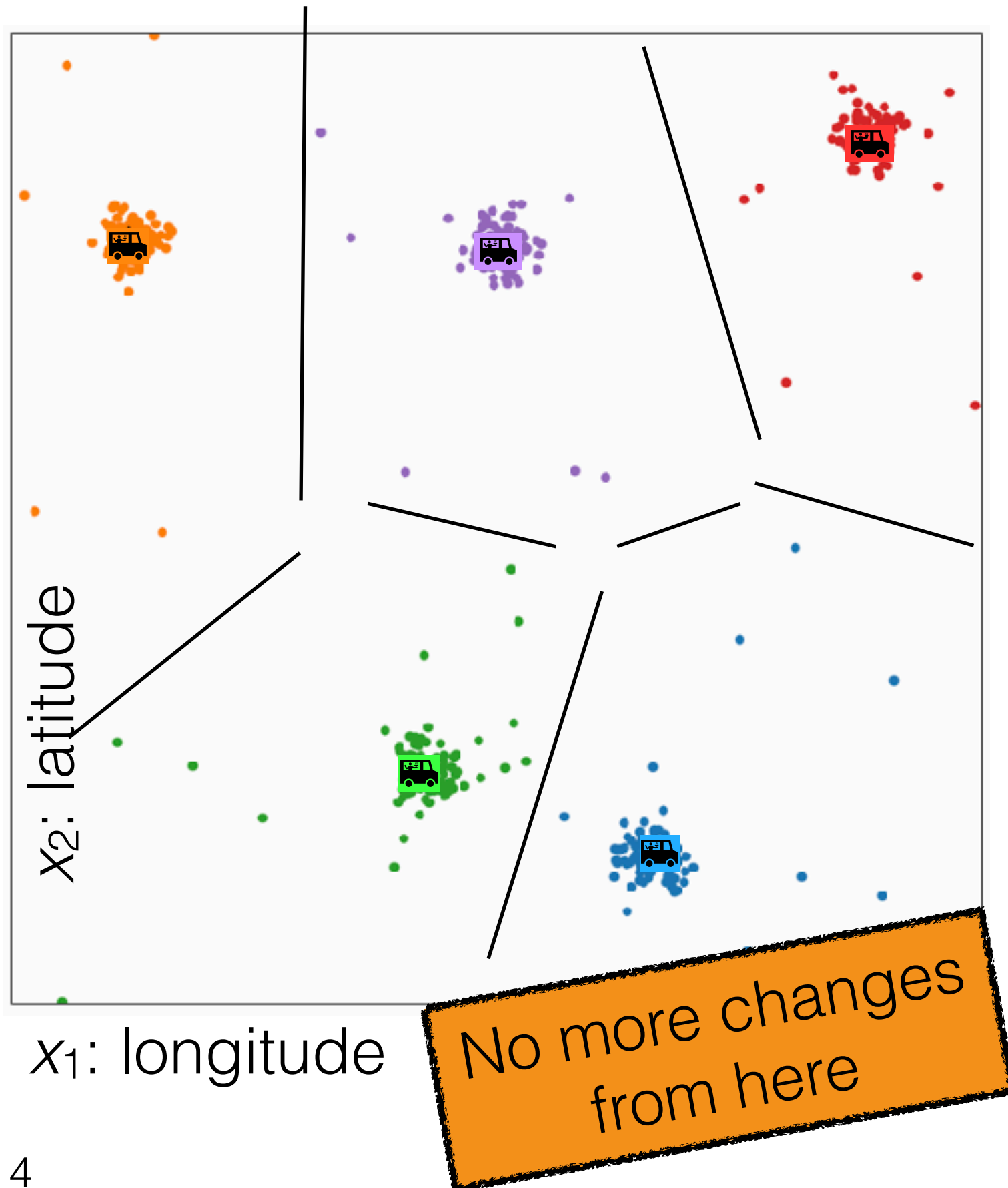$$\mu^{(j)} = \frac{\sum_{i=1}^n \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^n \mathbf{1}\{y^{(i)} = j\}}$$

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k, τ)
  Init {μ^(j)}_{j=1}^{k}
  for t = 1 to τ
```

$$\textbf{for } i = 1 \text{ to } n$$
$$y^{(i)} = \arg\min_j \| x^{(i)} - \mu^{(j)} \|_2^2$$
$$\textbf{for } j = 1 \text{ to } k$$
$$\mu^{(j)} = \frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

```
k-means(k,τ)
```
  Init $\{\mu^{(j)}\}_{j=1}^{k}$
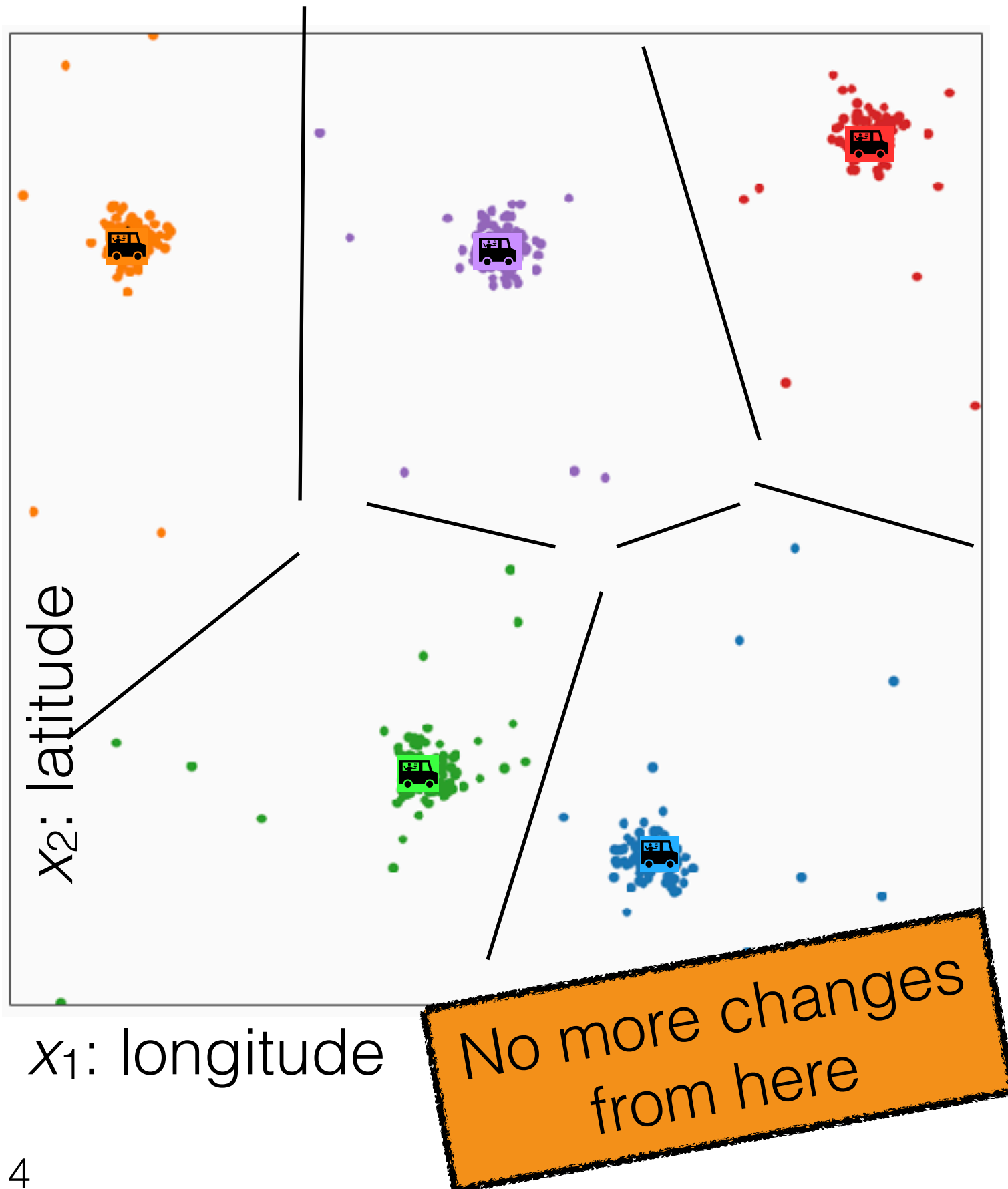  **for** t = 1 to $\tau$

    **for** i = 1 to n
      $y^{(i)} =$
        $\arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$
    **for** j = 1 to k
      $\mu^{(j)} =$
        $$\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

No more changes from here

```
k-means(k,τ)
```
$\text{Init } \{\mu^{(j)}\}_{j=1}^{k}$

**for** `t = 1 to` $\tau$
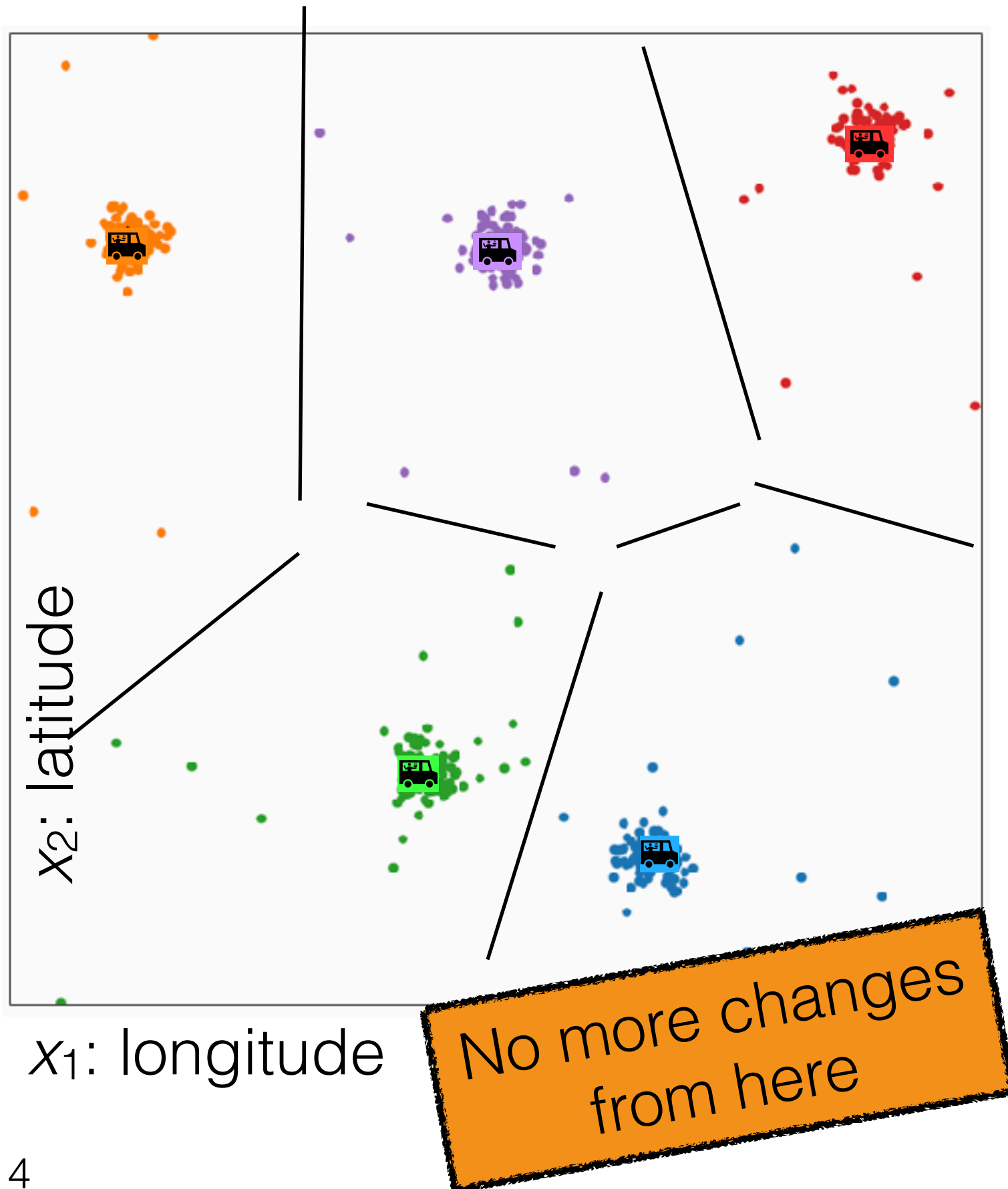
    **for** `i = 1 to n`

$$y^{(i)} = \arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$$

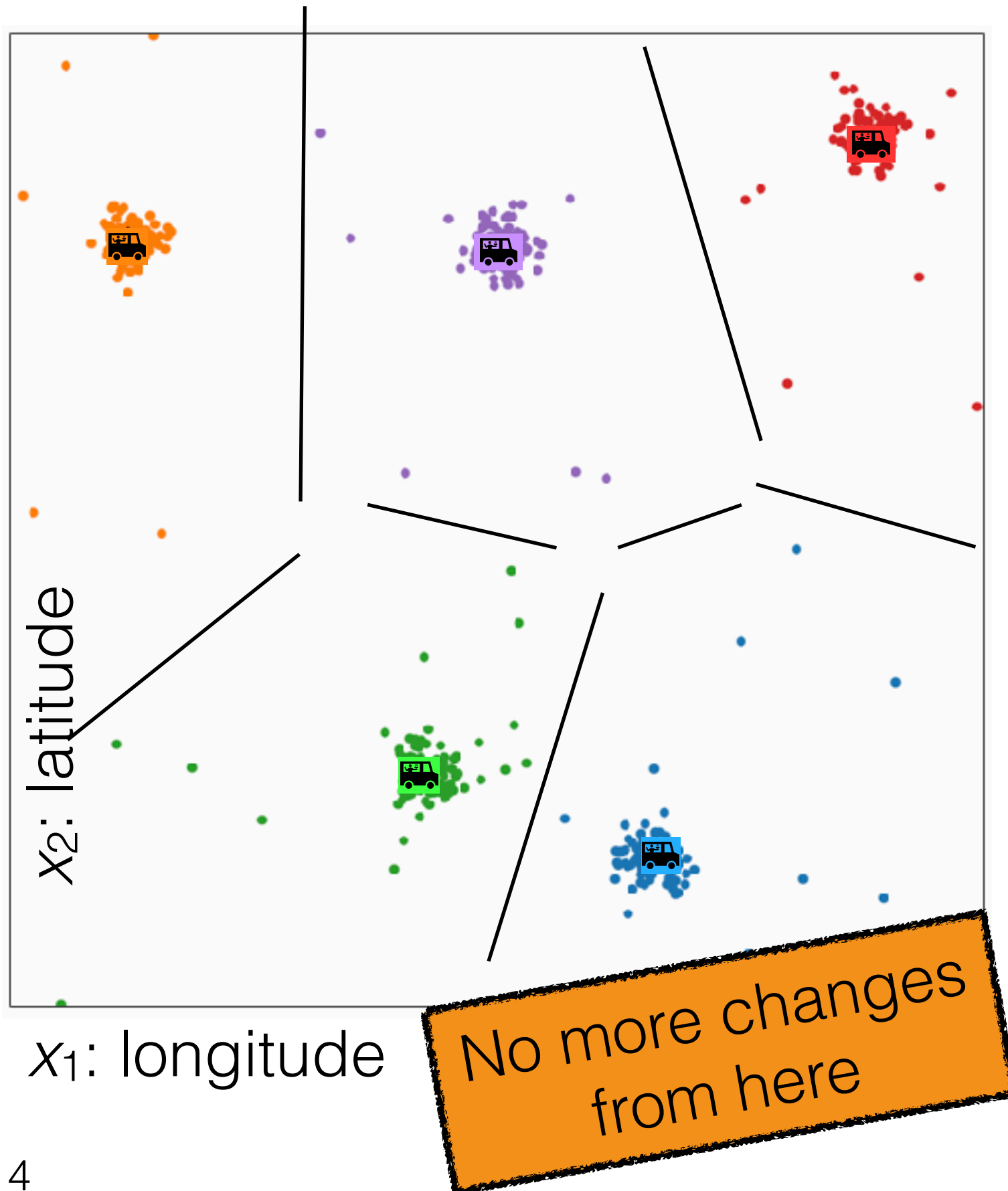    **for** `j = 1 to k`

$$\mu^{(j)} = \frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

No more changes from here

```
k-means(k,τ)
```
$\text{Init } \{\mu^{(j)}\}_{j=1}^{k}$

**for** t = 1 to $\tau$

    **for** i = 1 to n

$$y^{(i)} = \arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$$

    **for** j = 1 to k

$$\mu^{(j)} = \frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

How can I be so sure?

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

No more changes from here

```
k-means(k, τ)
```
Init $\{\mu^{(j)}\}_{j=1}^{k}$

**for** t = 1 to $\tau$

$y_{\text{old}} = y$

**for** i = 1 to n

$$y^{(i)} = \arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$$

**for** j = 1 to k

$$\mu^{(j)} = \frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$
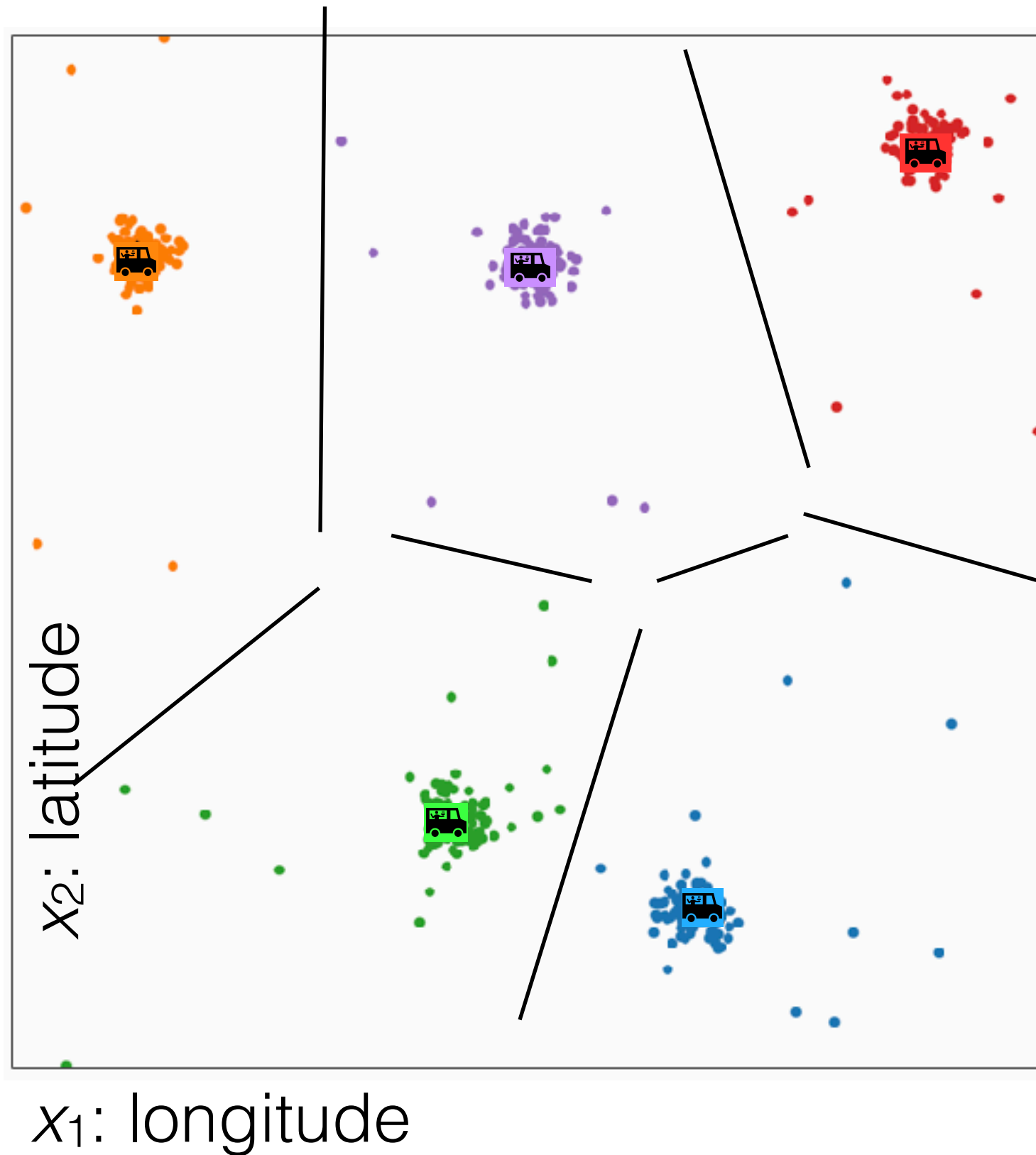
4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

No more changes from here

```
k-means(k,τ)
  Init {μ^(j)}_{j=1}^k
  for t = 1 to τ
    y_old = y
    for i = 1 to n
      y^(i) =
        arg min_j ||x^(i) - μ^(j)||_2^2
    for j = 1 to k
      μ^(j) =
        (Σ_{i=1}^n 1{y^(i) = j}x^(i)) / (Σ_{i=1}^n 1{y^(i) = j})
    if y = y_old
      break
```

$$\text{k-means}(\texttt{k},\tau)$$
$$\text{Init } \{\mu^{(j)}\}_{j=1}^k$$
$$\textbf{for } \texttt{t = 1 to } \tau$$
$$y_{\text{old}} = y$$
$$\textbf{for } \texttt{i = 1 to n}$$
$$y^{(i)} = \arg\min_j \|x^{(i)} - \mu^{(j)}\|_2^2$$
$$\textbf{for } \texttt{j = 1 to k}$$
$$\mu^{(j)} = \frac{\sum_{i=1}^n \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^n \mathbf{1}\{y^{(i)} = j\}}$$
$$\textbf{if } y = y_{\text{old}}$$
$$\texttt{break}$$

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

No more changes from here

```
k-means(k,τ)
  Init {μ^(j)}^k_{j=1}, {y^(i)}^n_{i=1}
  for t = 1 to τ
      y_old = y
      for i = 1 to n
      y^(i) =
          arg min_j ‖x^(i) − μ^(j)‖^2_2
      for j = 1 to k
      μ^(j) =
          (∑^n_{i=1} 1{y^(i) = j}x^(i)) / (∑^n_{i=1} 1{y^(i) = j})
  if y = y_old
      break
```

4

# k-means algorithm



$x_2$: latitude

$x_1$: longitude

No more changes from here

```
k-means(k,τ)
```

$\text{Init } \{\mu^{(j)}\}_{j=1}^{k}, \{y^{(i)}\}_{i=1}^{n}$

**for** t = 1 to τ

$\quad y_{\text{old}} = y$

**for** i = 1 to n

$\quad y^{(i)} =$
$\quad \arg\min_{j} \|x^{(i)} - \mu^{(j)}\|_2^2$

**for** j = 1 to k

$\quad \mu^{(j)} =$
$$\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$$

**if** $y = y_{\text{old}}$

break

**return** $\{\mu^{(j)}\}_{j=1}^{k}, \{y^{(i)}\}_{i=1}^{n}$

4

# Compare to classification



$x_2$: latitude

$x_1$: longitude

5

# Compare to classification



- Did we just do *k*-class classification?

$x_2$: latitude

$x_1$: longitude

# Compare to classification
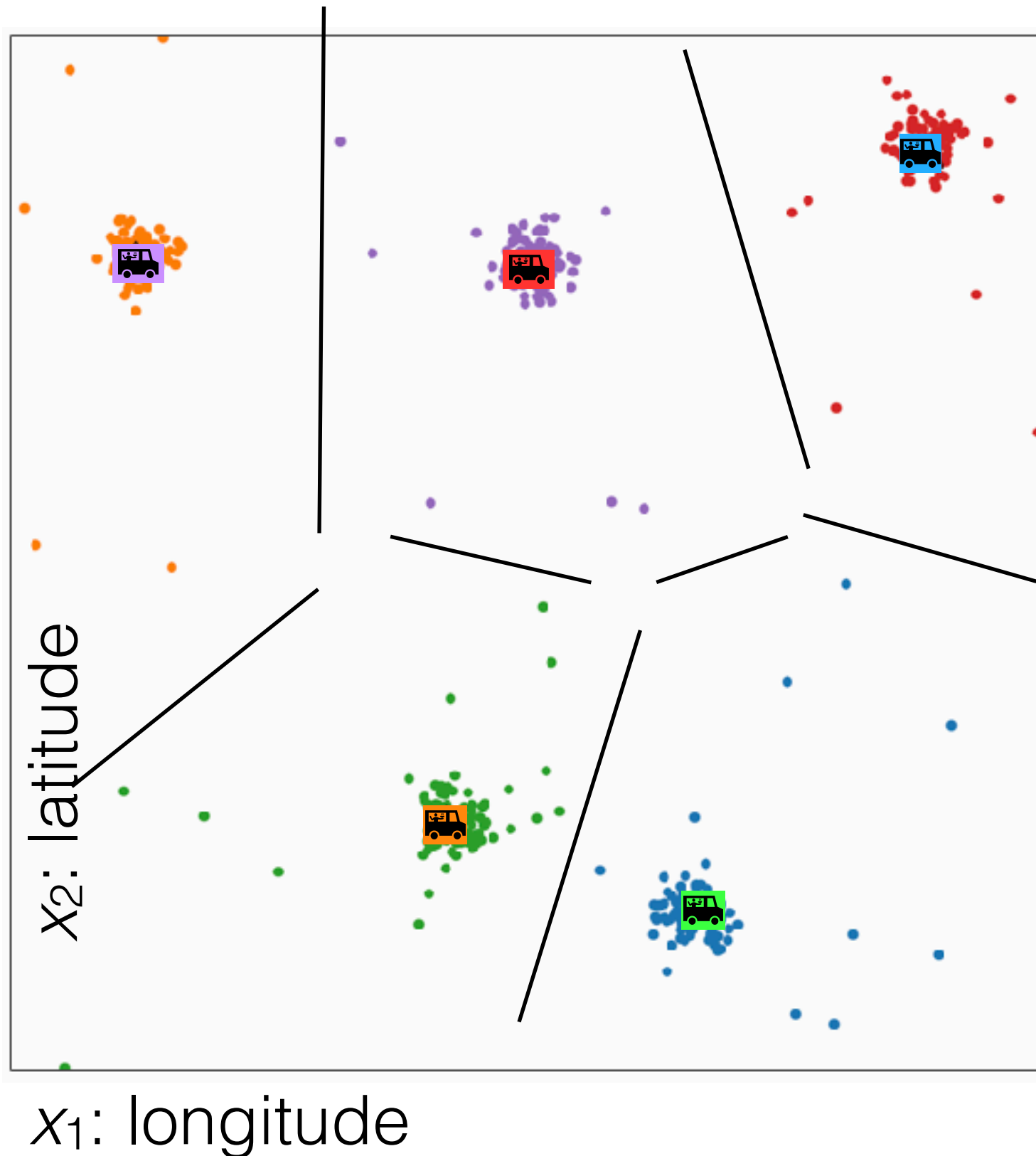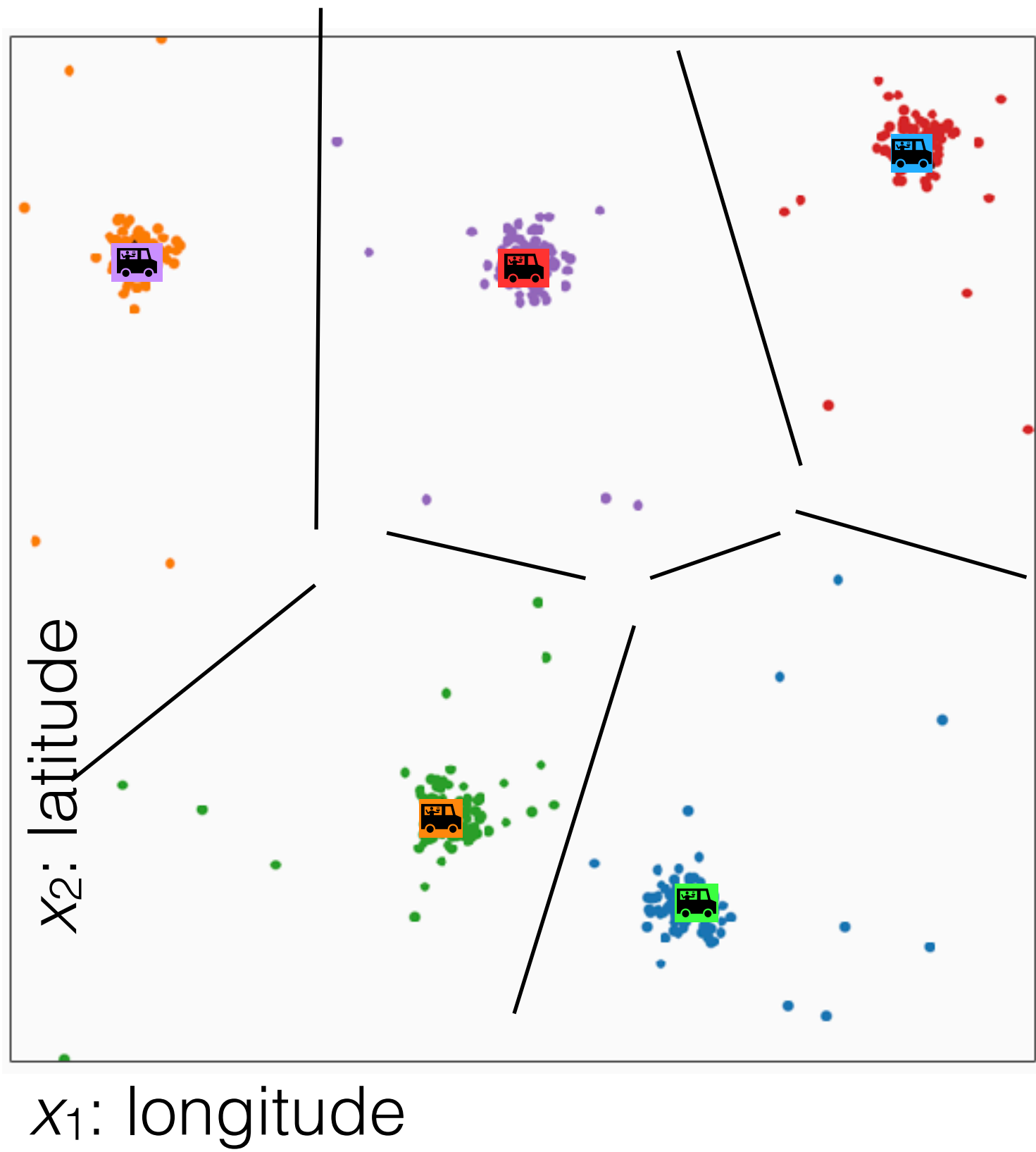


$x_2$: latitude

$x_1$: longitude

- Did we just do *k*-class classification?
- Looks like we assigned a label $y^{(i)}$, which takes *k* different values, to each feature vector $x^{(i)}$
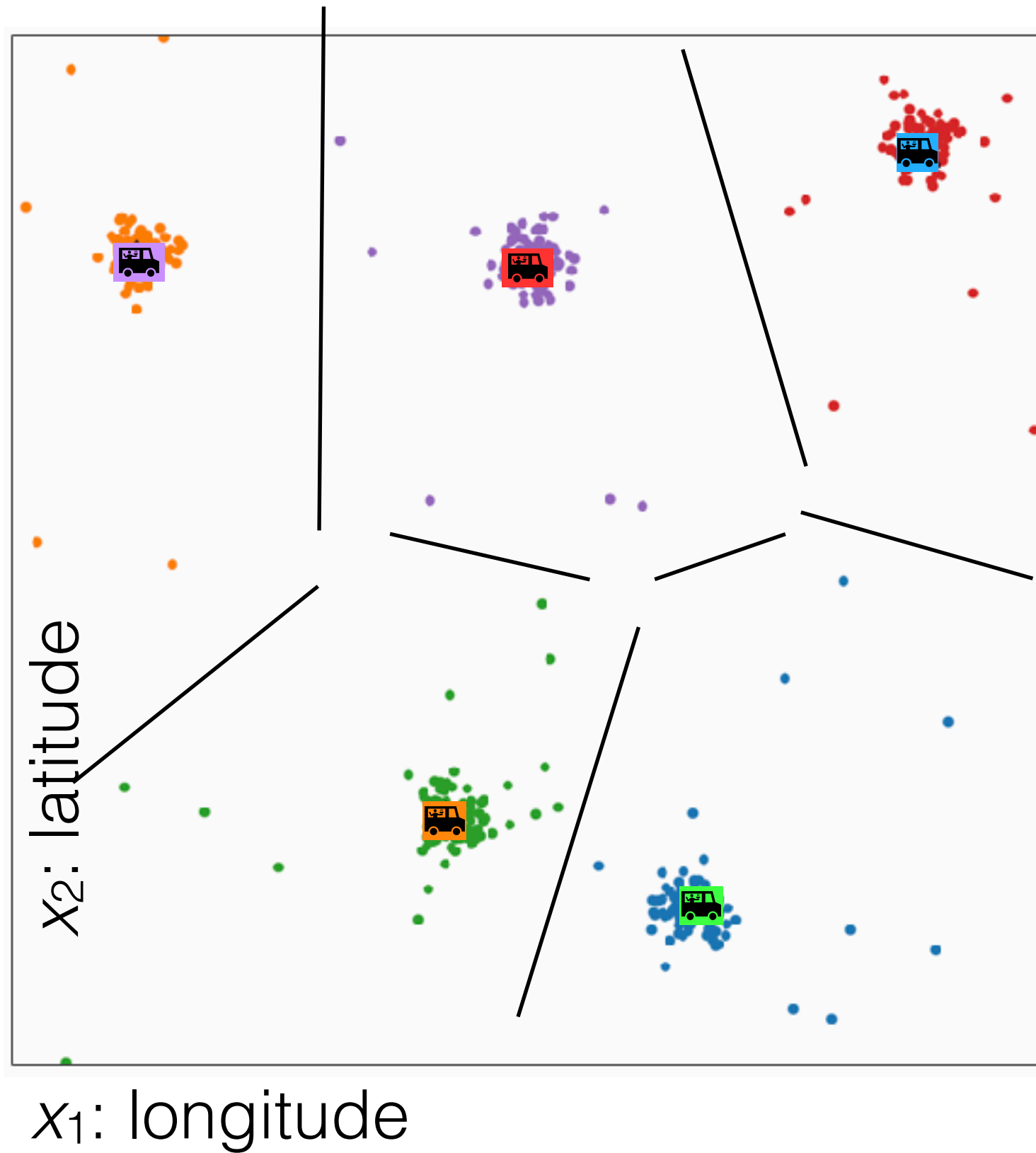
# Compare to classification
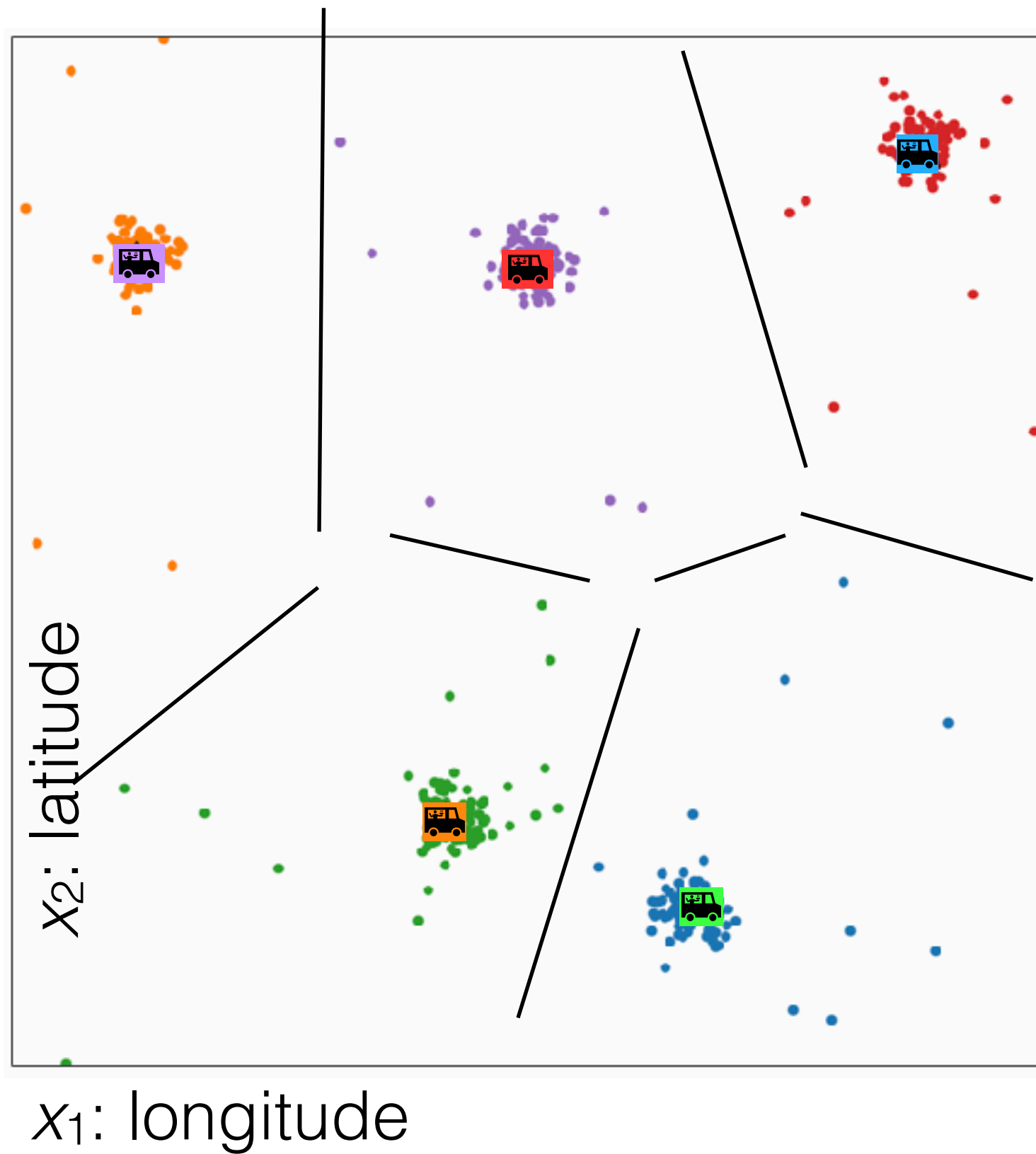


$x_2$: latitude

$x_1$: longitude

- Did we just do *k*-class classification?
- Looks like we assigned a label $y^{(i)}$, which takes *k* different values, to each feature vector $x^{(i)}$
- But we didn't use any labeled data

# Compare to classification



$x_2$: latitude

$x_1$: longitude

- Did we just do *k*-class classification?
- Looks like we assigned a label $y^{(i)}$, which takes *k* different values, to each feature vector $x^{(i)}$
- But we didn't use any labeled data

# Compare to classification

- Did we just do *k*-class classification?

- Looks like we assigned a label $y^{(i)}$, which takes *k* different values, to each feature vector $x^{(i)}$

- But we didn't use any labeled data

- The "labels" here don't have meaning; I could permute them and have the same result

$x_2$: latitude

$x_1$: longitude

# Compare to classification



$x_2$: latitude

$x_1$: longitude

- Did we just do *k*-class classification?

- Looks like we assigned a label $y^{(i)}$, which takes *k* different values, to each feature vector $x^{(i)}$

- But we didn't use any labeled data

- The "labels" here don't have meaning; I could permute them and have the same result

# Compare to classification



x₂: latitude

x₁: longitude

- Did we just do *k*-class classification?

- Looks like we assigned a label $y^{(i)}$, which takes *k* different values, to each feature vector $x^{(i)}$

- But we didn't use any labeled data

- The "labels" here don't have meaning; I could permute them and have the same result

# Compare to classification



$x_2$: latitude

$x_1$: longitude

- Did we just do *k*-class classification?

- Looks like we assigned a label $y^{(i)}$, which takes *k* different values, to each feature vector $x^{(i)}$

- But we didn't use any labeled data

- The "labels" here don't have meaning; I could permute them and have the same result

- Output is really a *partition* of the data

$x_2$: latitude

$x_1$: longitude

6

- So what did we do?

$x_2$: latitude

$x_1$: longitude

6

- So what did we do?
- We *clustered* the data

$x_2$: latitude

$x_1$: longitude

6

- So what did we do?
- We *clustered* the data: we grouped the data by similarity

$x_2$: latitude

$x_1$: longitude

$x_2$: latitude

$x_1$: longitude

- So what did we do?
- We *clustered* the data: we grouped the **data** by similarity
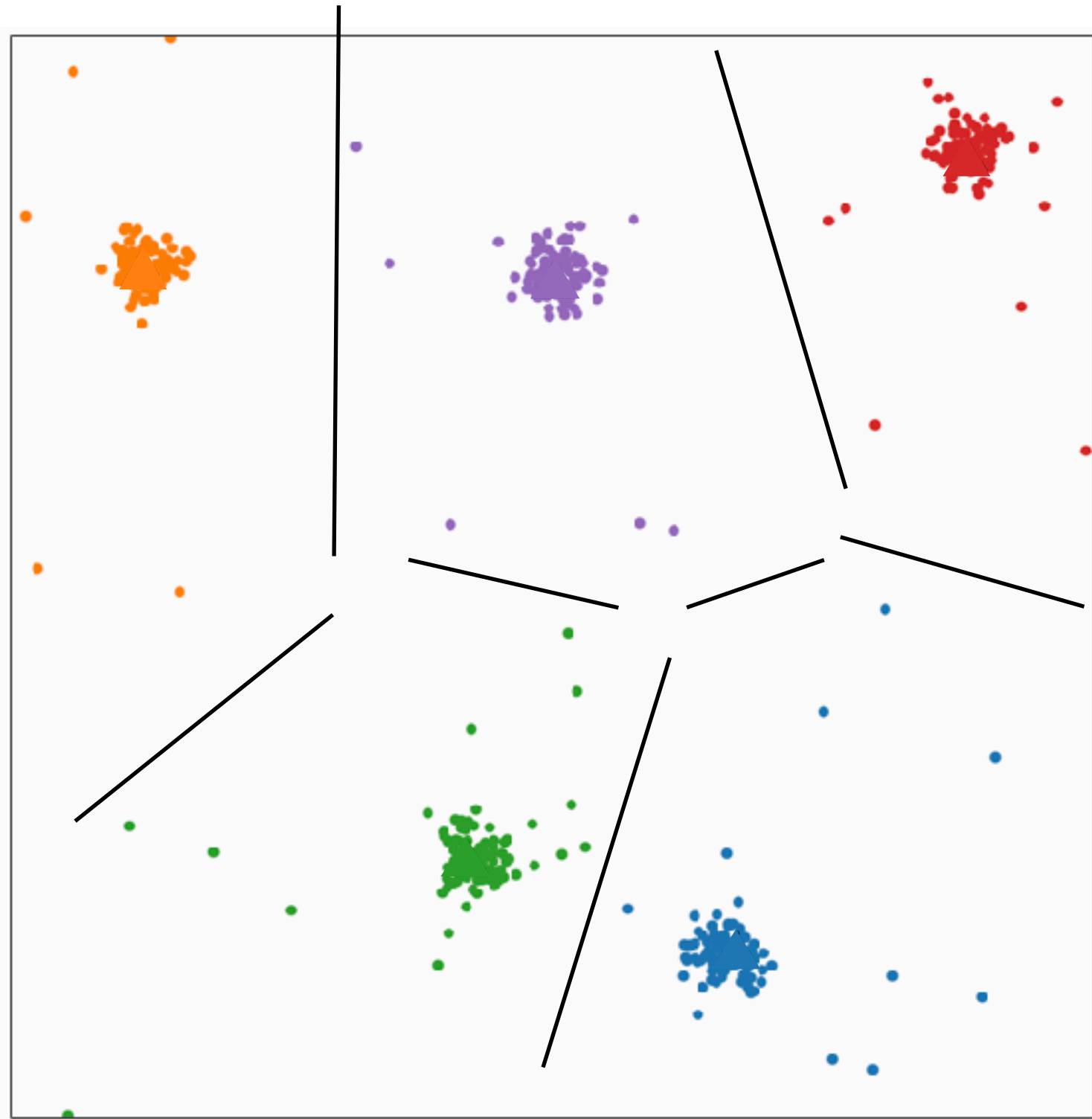
6

$x_2$: latitude

$x_1$: longitude

- So what did we do?
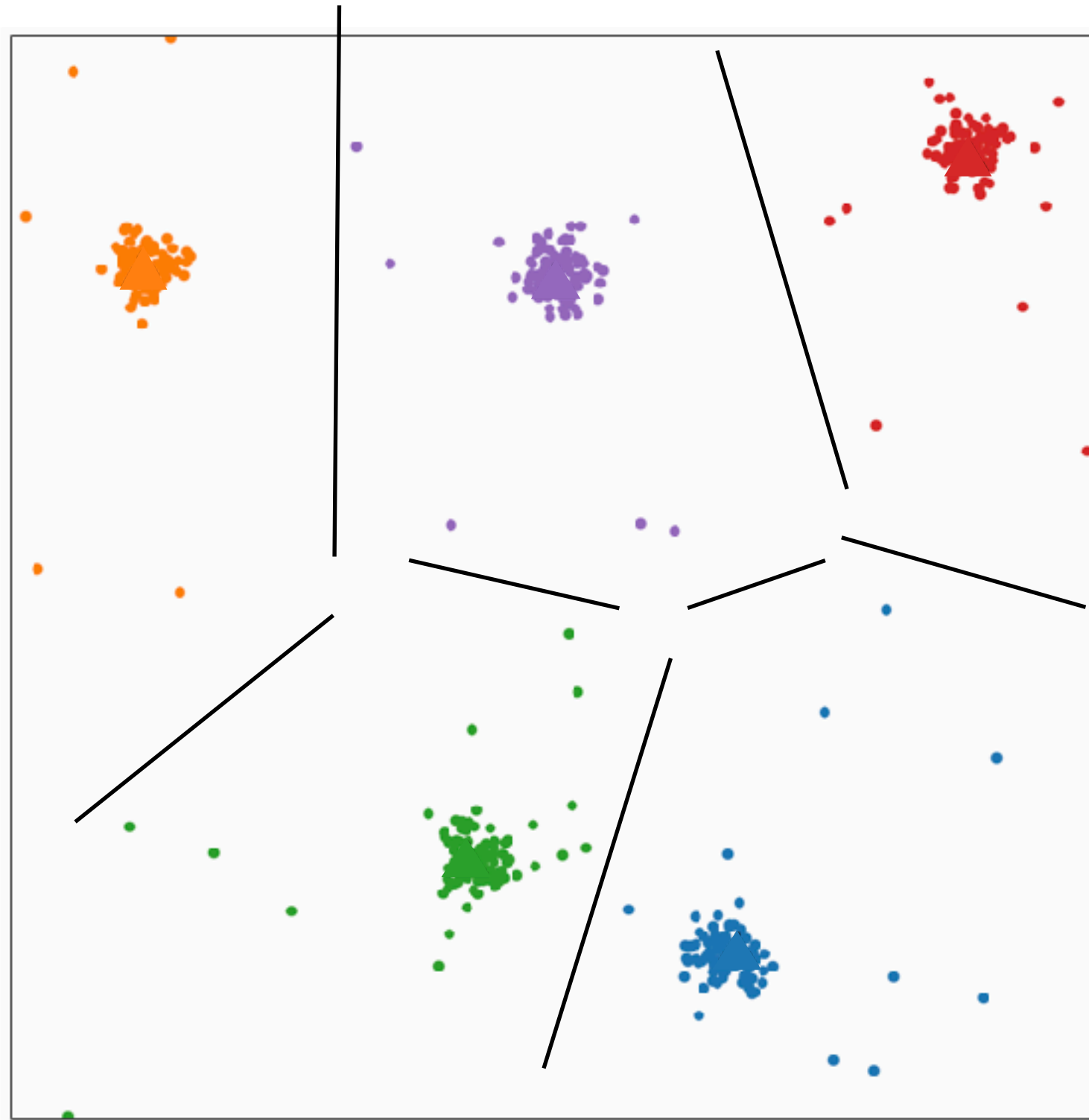- We *clustered* the data: we grouped the **data** by similarity

- So what did we do?
- We *clustered* the data: we grouped the **data** by similarity

- So what did we do?
- We *clustered* the data: we grouped the data by **similarity**

- So what did we do?
- We *clustered* the data: we **grouped** the data by similarity

- So what did we do?
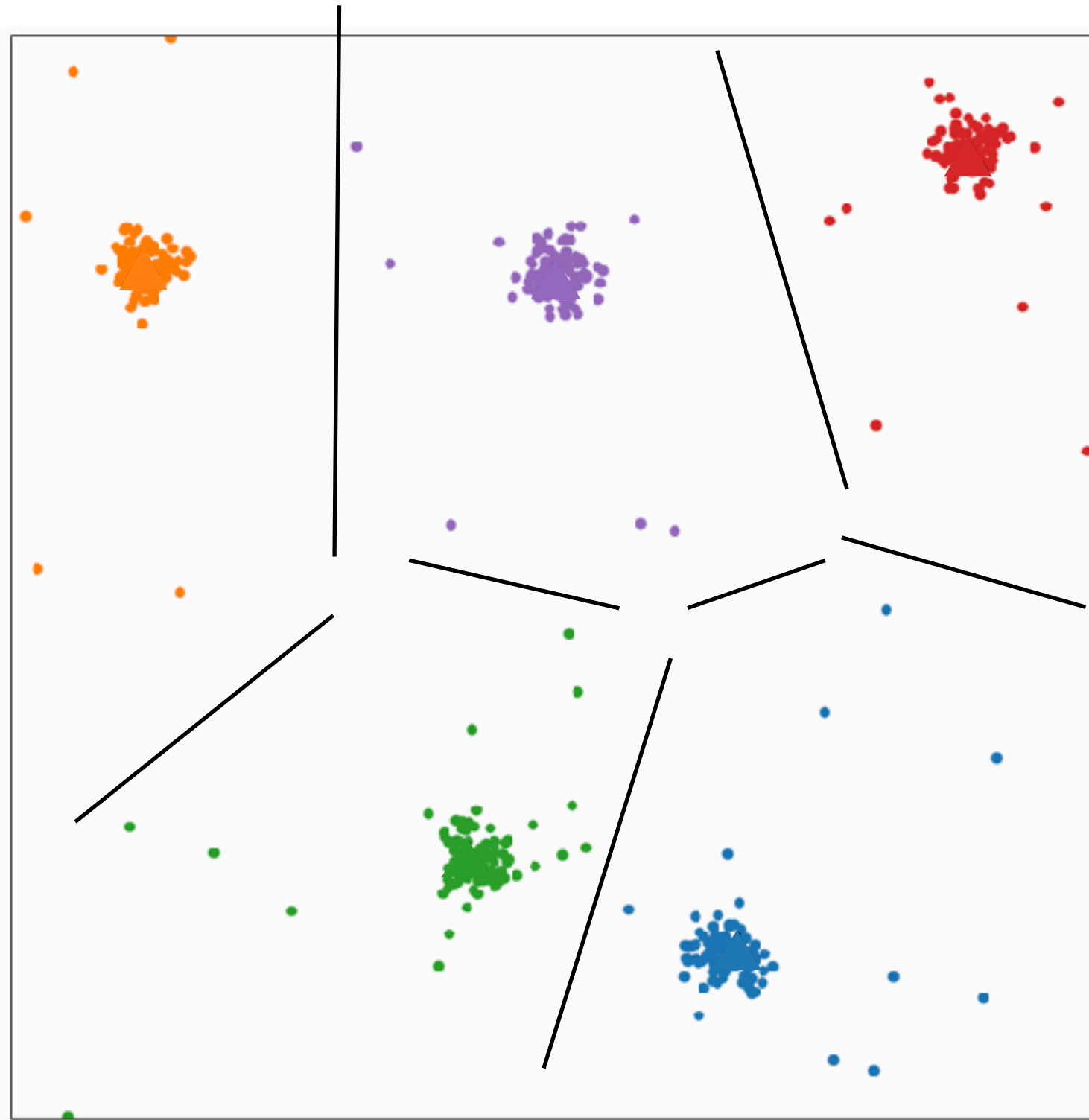- We *clustered* the data: we **grouped** the data by similarity

- So what did we do?
- We *clustered* the data: we **grouped** the data by similarity

- So what did we do?
- We *clustered* the data: we **grouped** the data by similarity

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
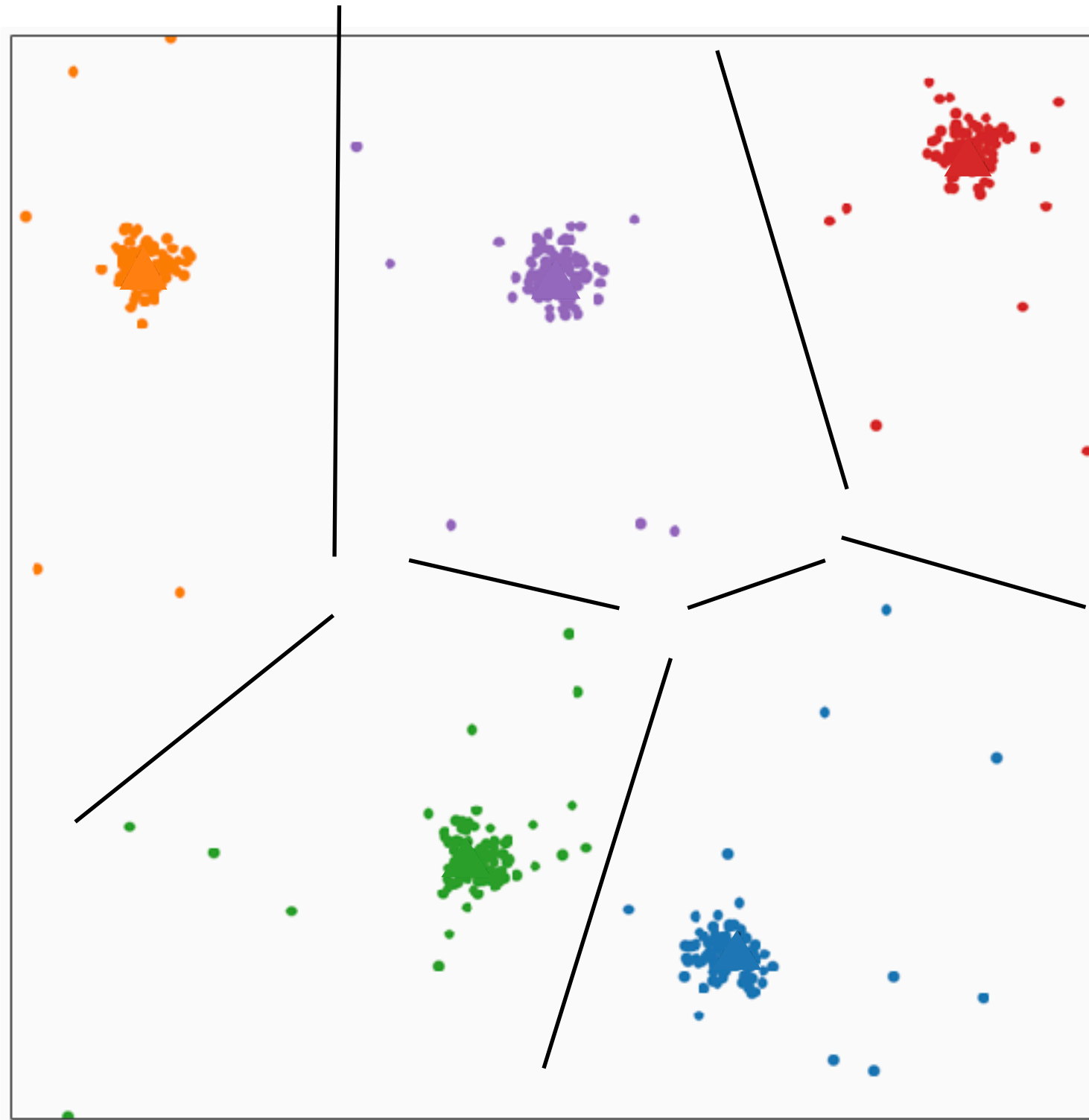  - Why not just plot the data?

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
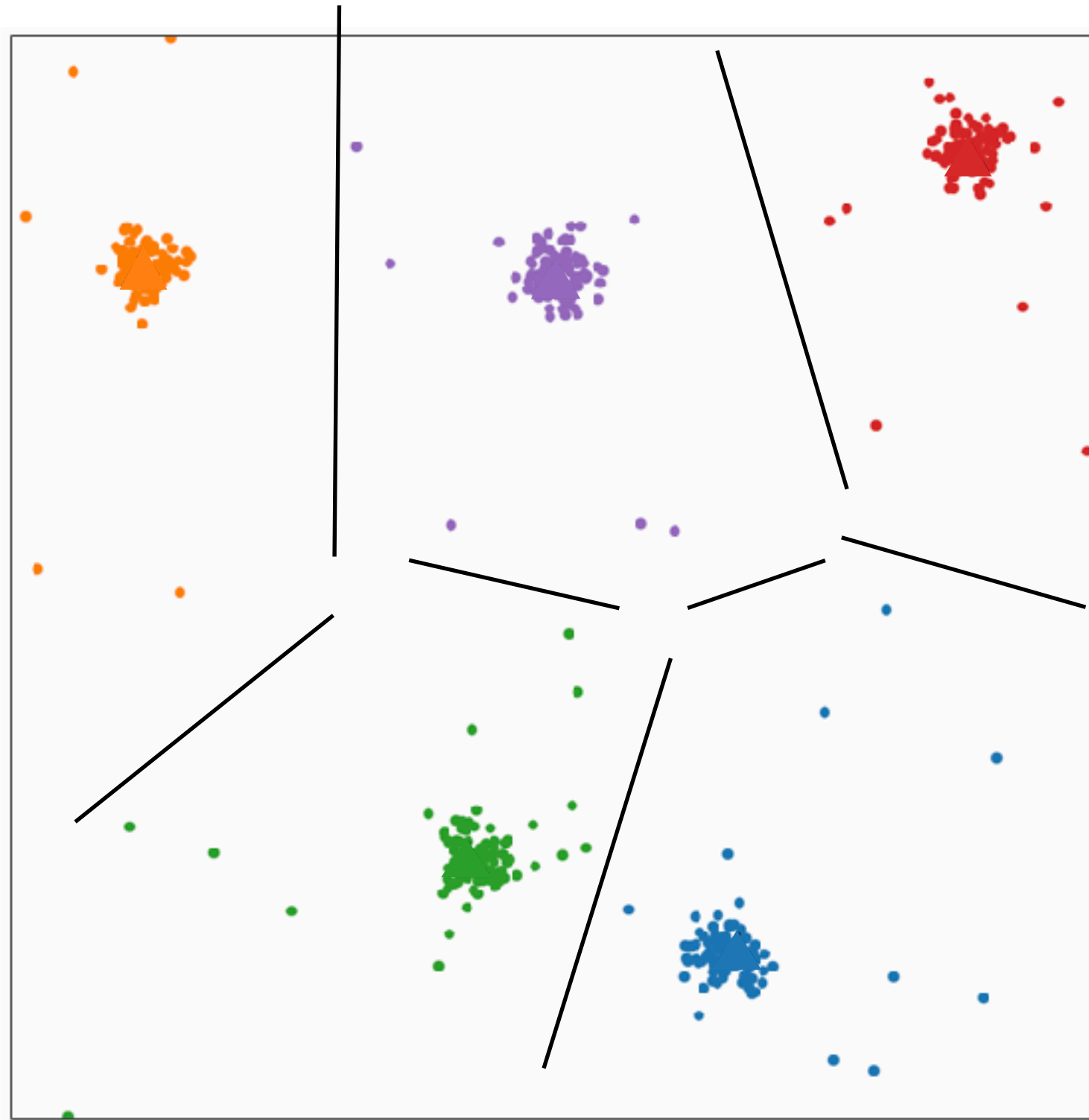  - Why not just plot the data? You should!

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also:
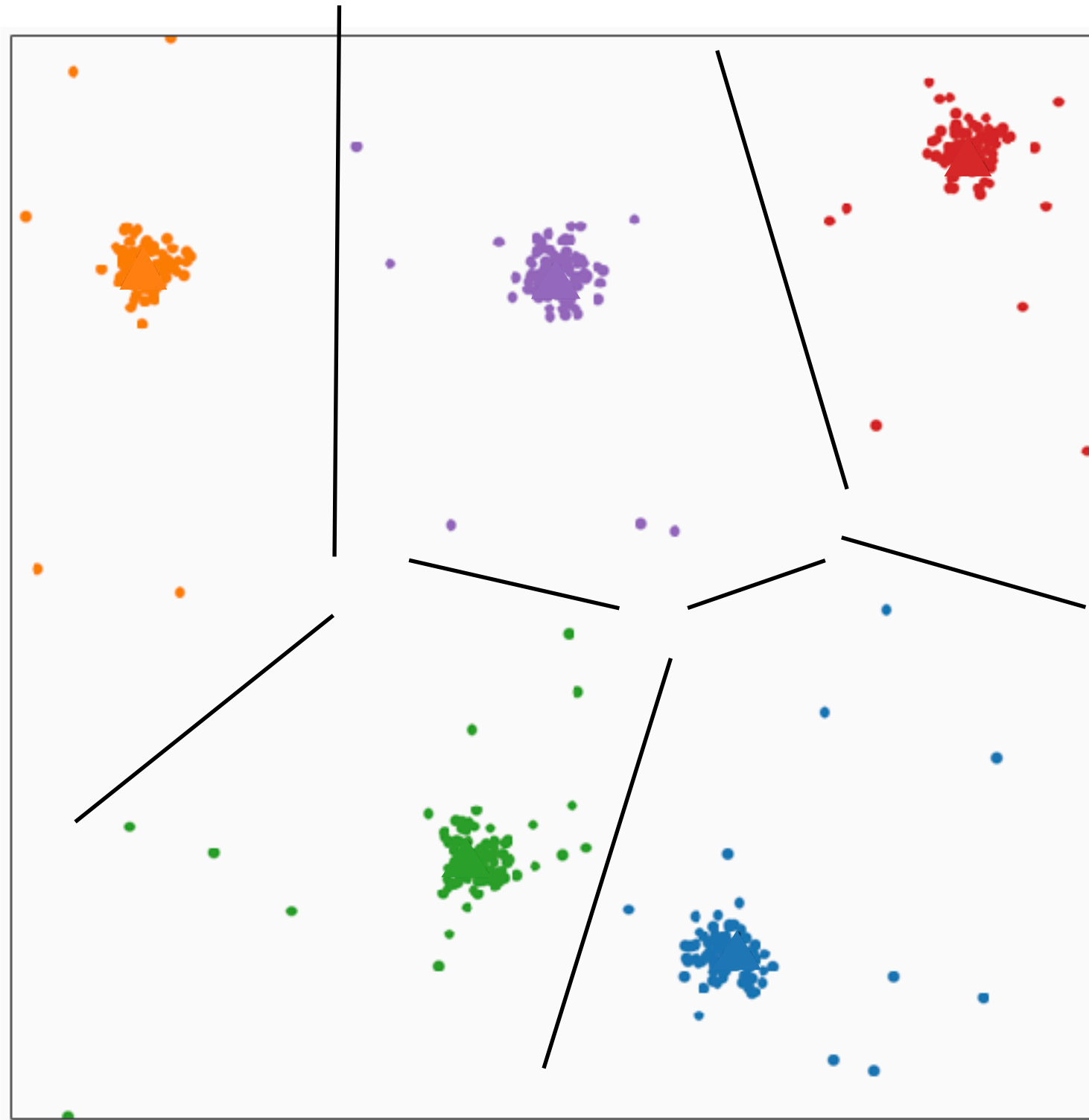
- So what did we do?
- We *clustered* the data: we grouped the data by similarity
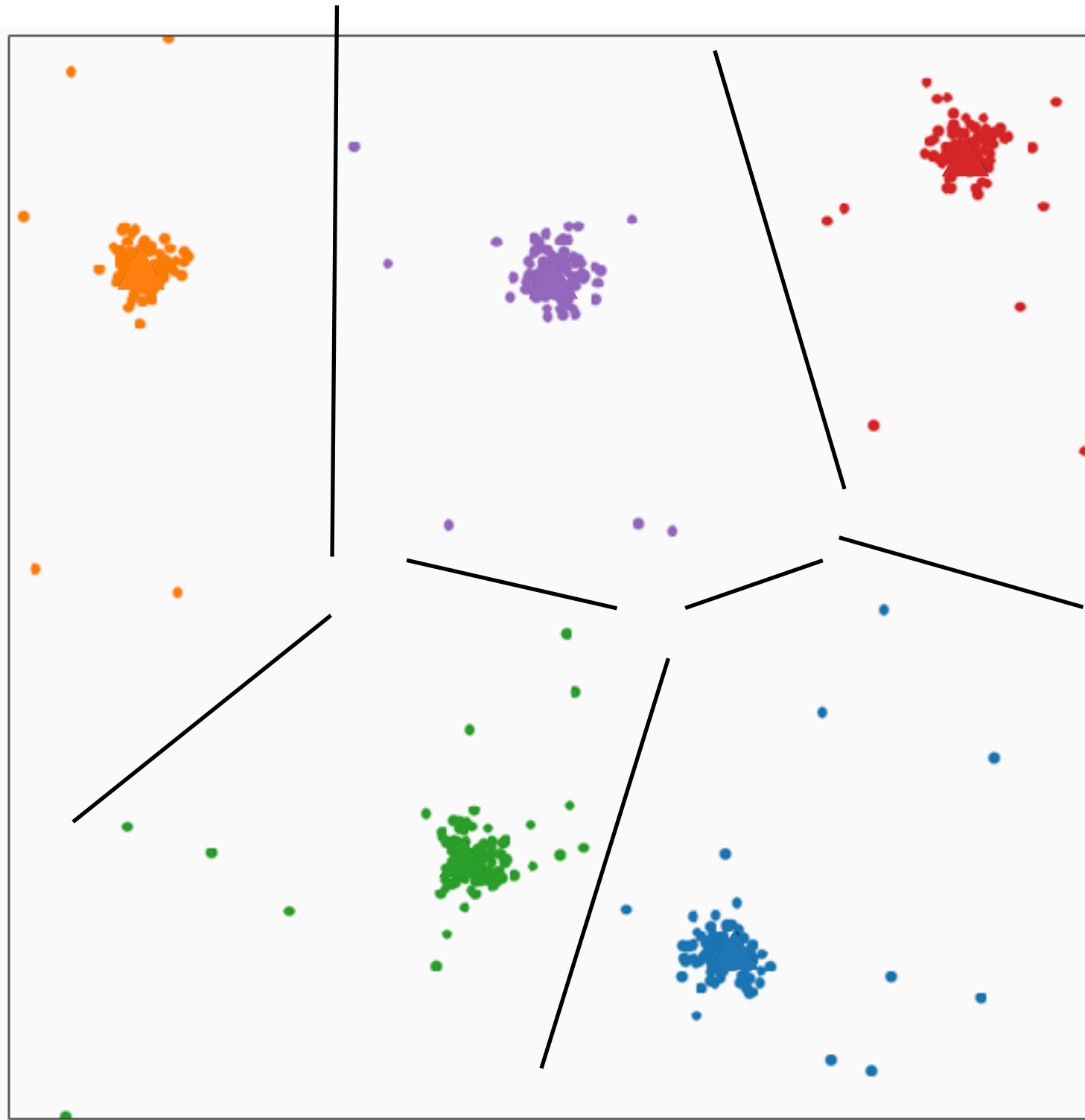  - Why not just plot the data? You should! But also: Precision

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also: Precision, big data

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also: Precision, big data, high dimensions

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also: Precision, big data, high dimensions, high volume

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also: Precision, big data, high dimensions, high volume
- An example of *unsupervised learning*: no labeled data, & we're finding patterns

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also: Precision, big data, high dimensions, high volume
- An example of *unsupervised learning*: no labeled data, & we're finding patterns
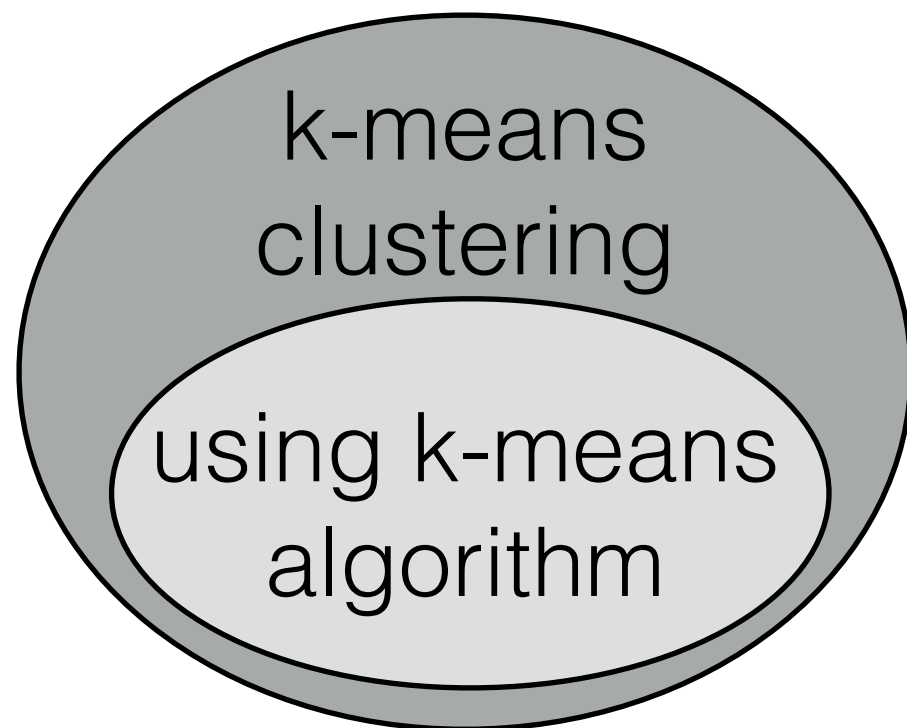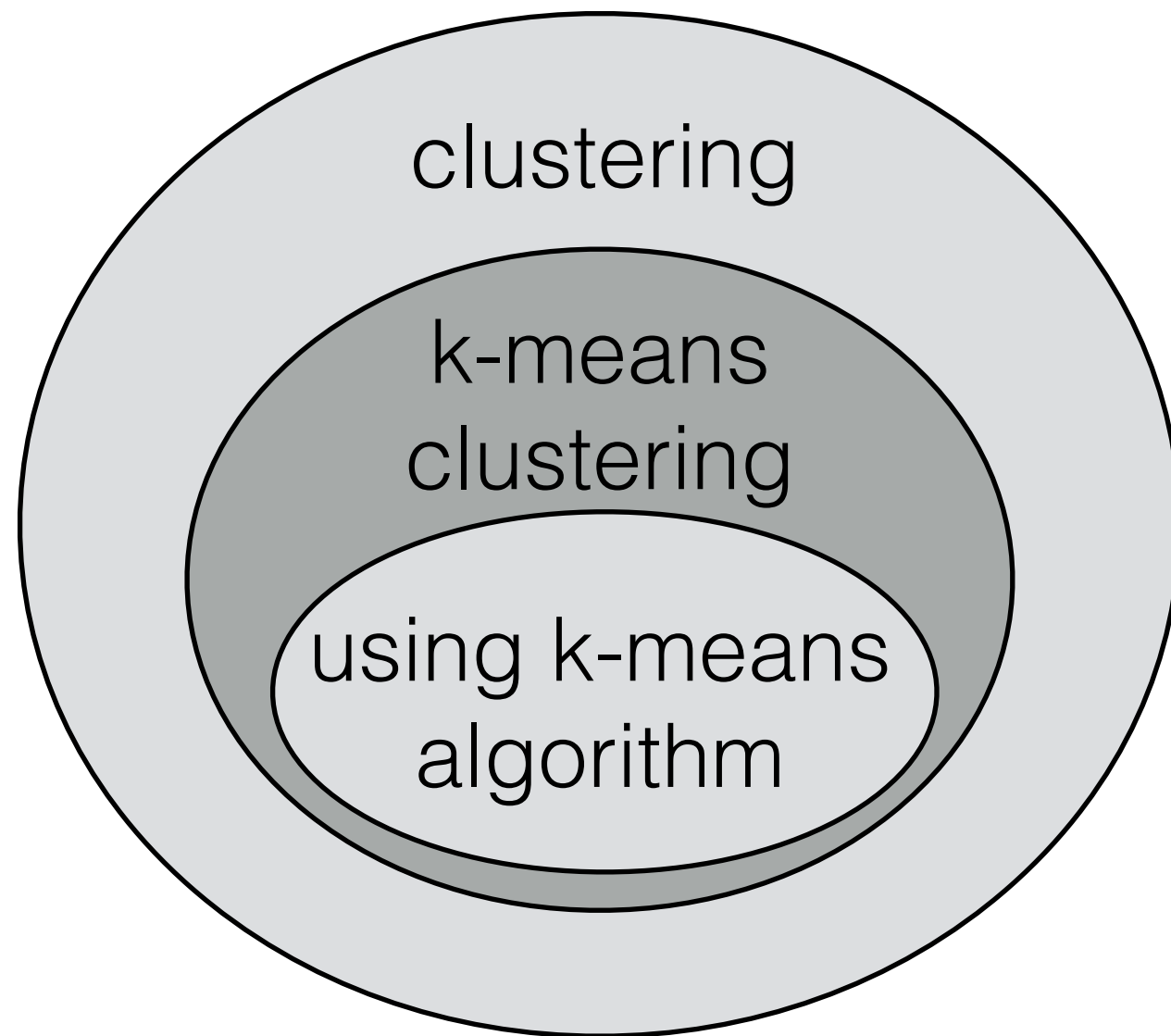
# Clustering & related

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also: Precision, big data, high dimensions, high volume
- An example of *unsupervised learning*: no labeled data, & we're finding patterns
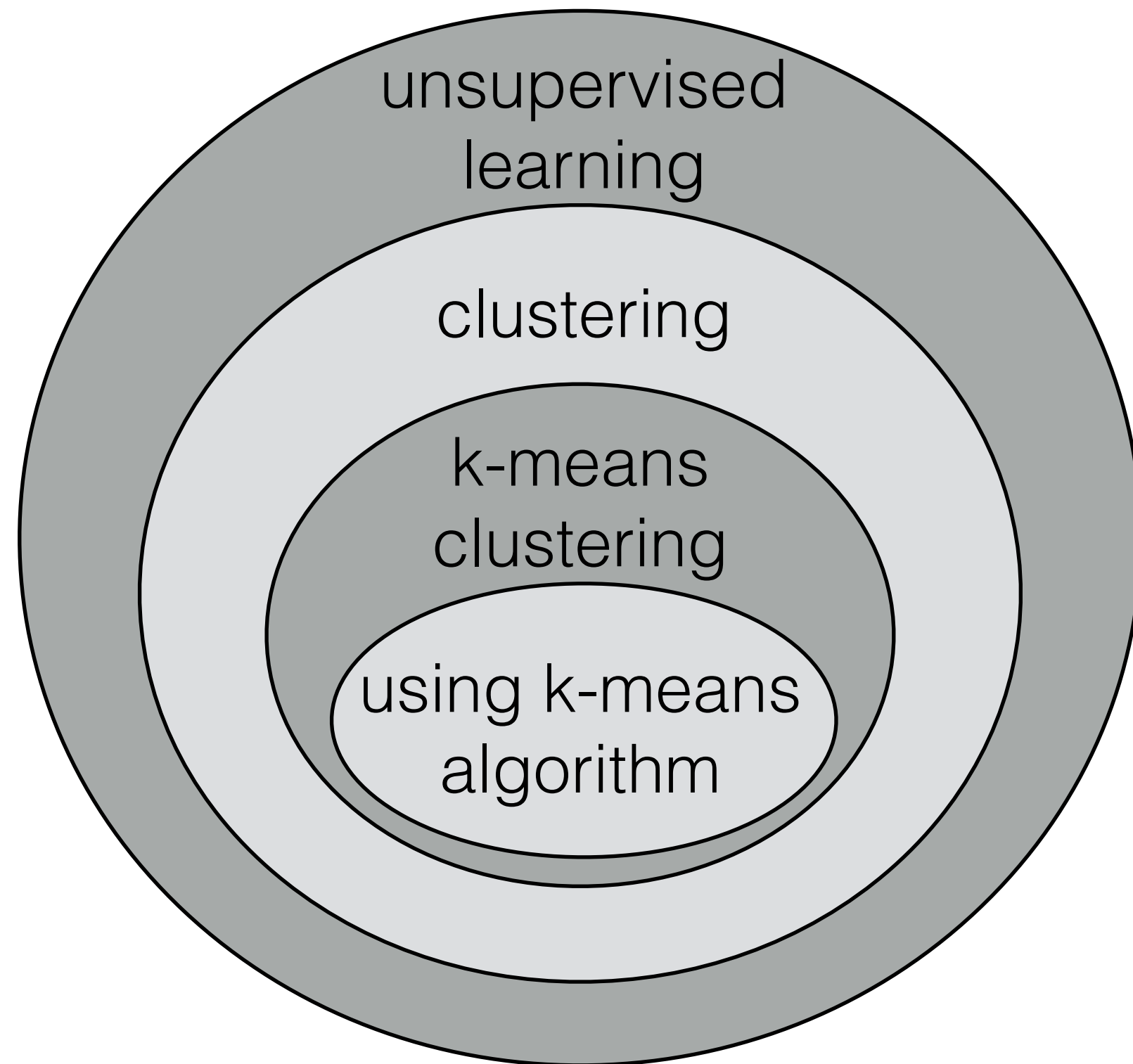
# Clustering & related

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also: Precision, big data, high dimensions, high volume
- An example of *unsupervised learning*: no labeled data, & we're finding patterns

using k-means algorithm
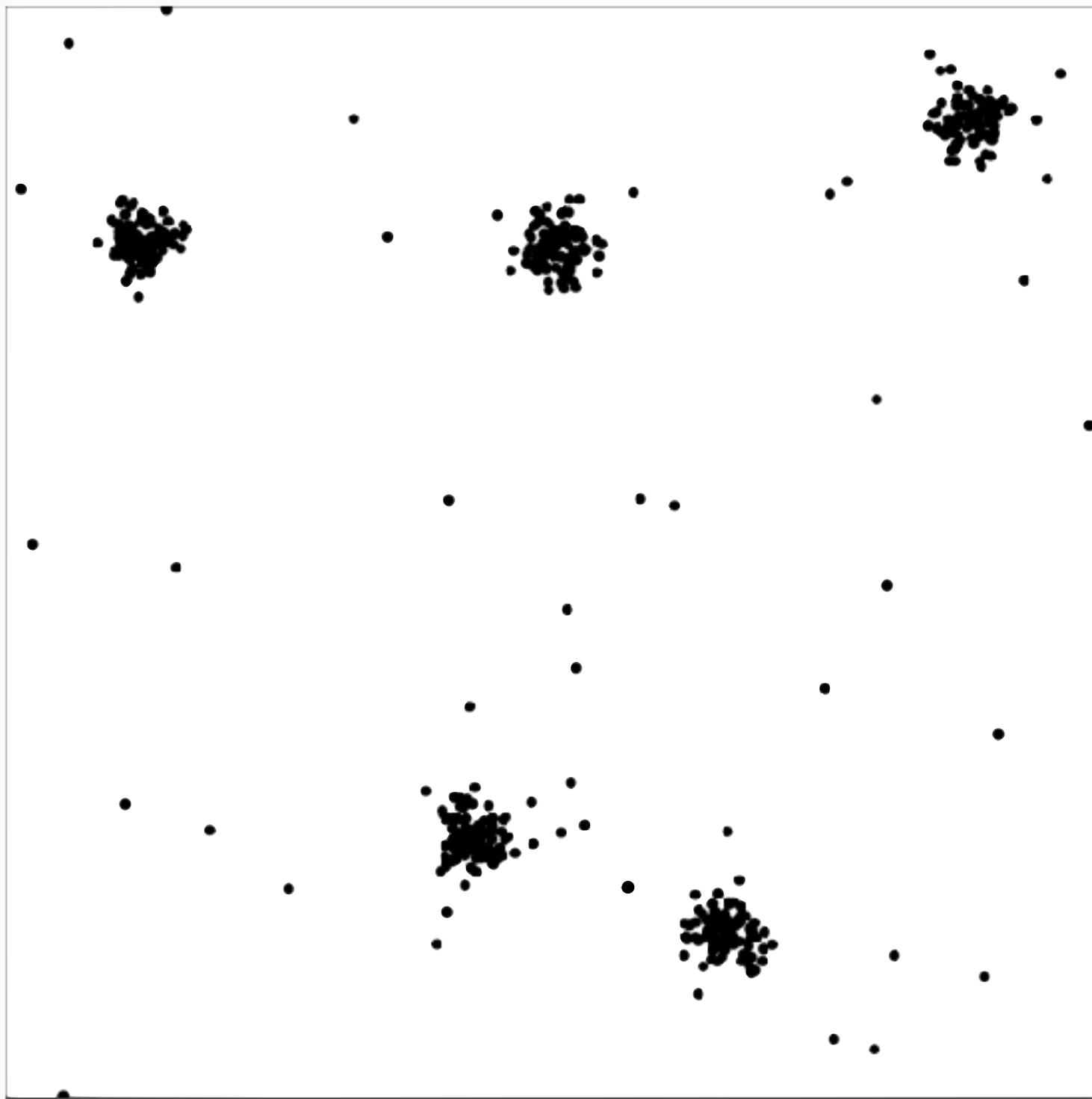
# Clustering & related

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also: Precision, big data, high dimensions, high volume
- An example of *unsupervised learning*: no labeled data, & we're finding patterns

k-means clustering

using k-means algorithm

# Clustering & related



clustering

k-means clustering

using k-means algorithm

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also: Precision, big data, high dimensions, high volume
- An example of *unsupervised learning*: no labeled data, & we're finding patterns

# Clustering & related

unsupervised learning

clustering

k-means clustering

using k-means algorithm

- So what did we do?
- We *clustered* the data: we grouped the data by similarity
  - Why not just plot the data? You should! But also: Precision, big data, high dimensions, high volume
- An example of *unsupervised learning*: no labeled data, & we're finding patterns

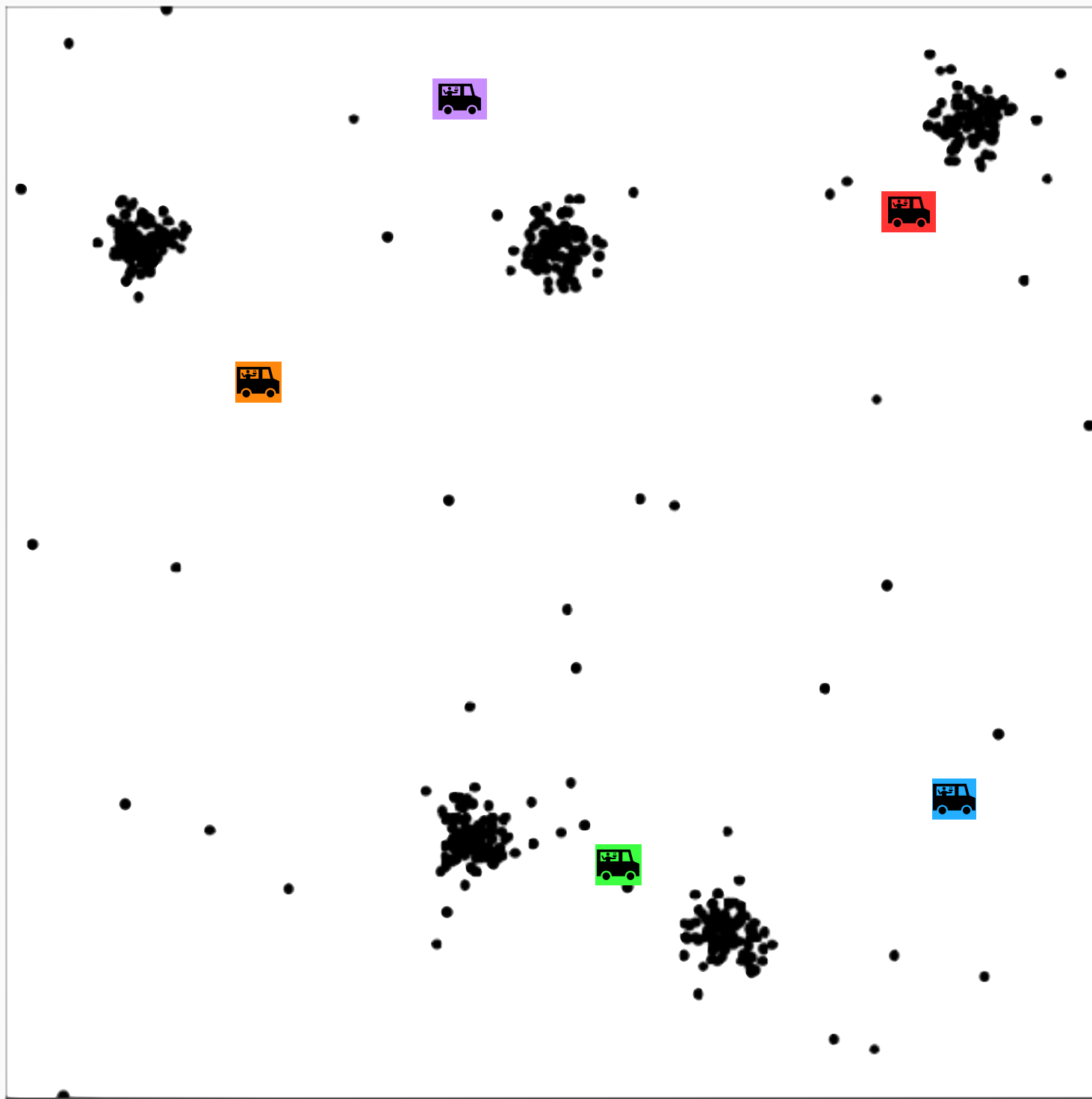# k-means algorithm: initialization

# k-means algorithm: initialization

- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective
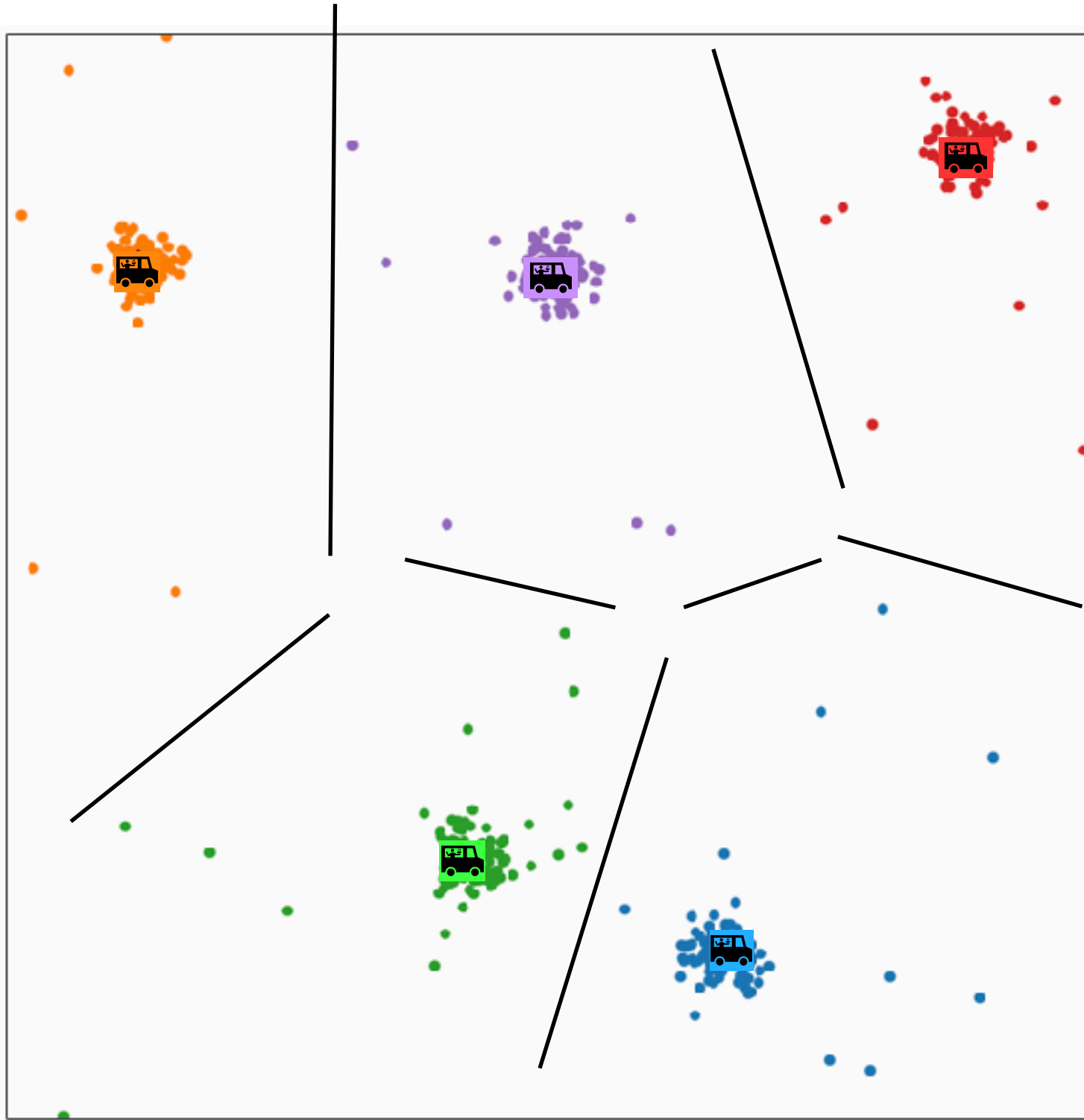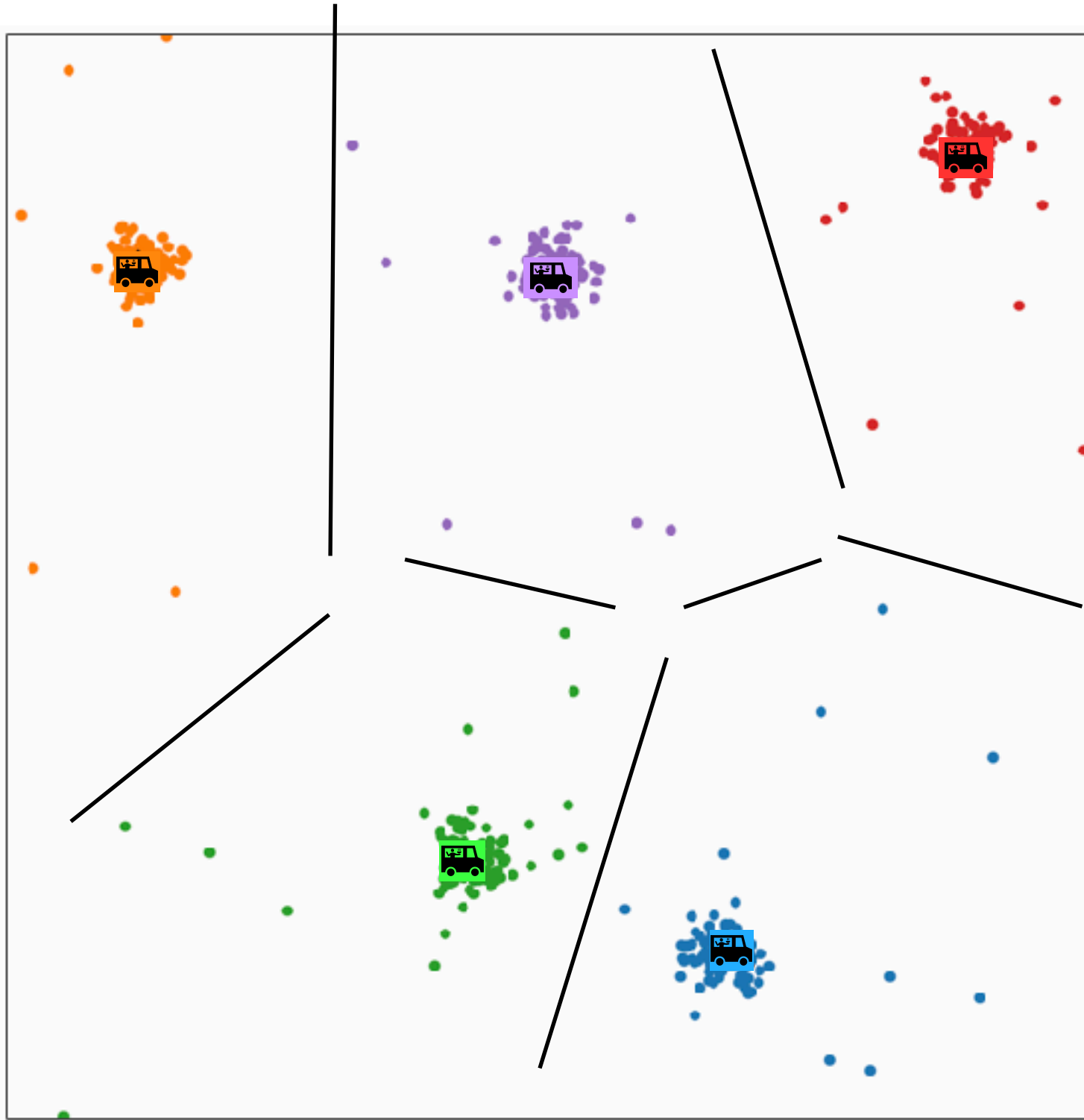
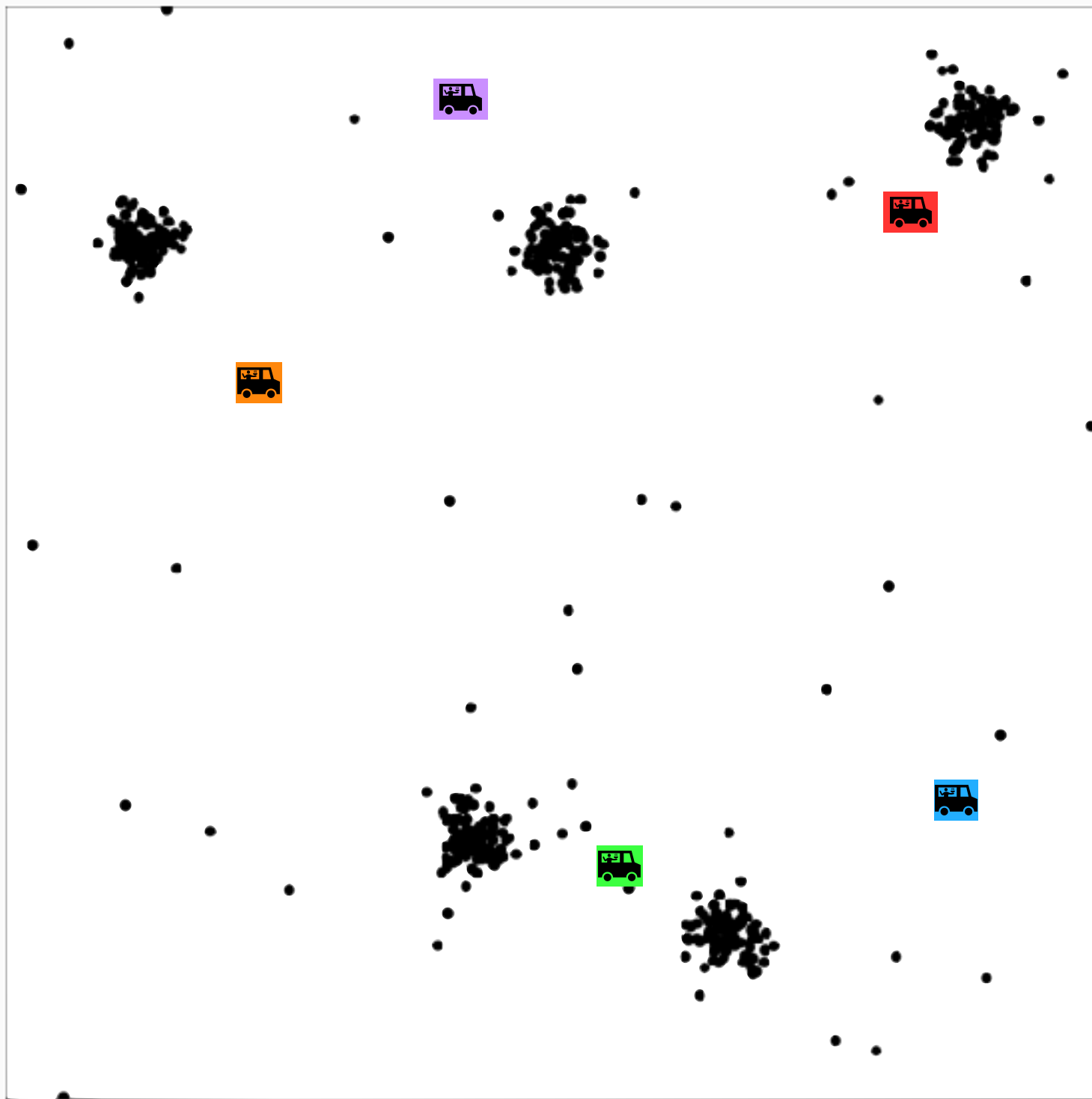# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

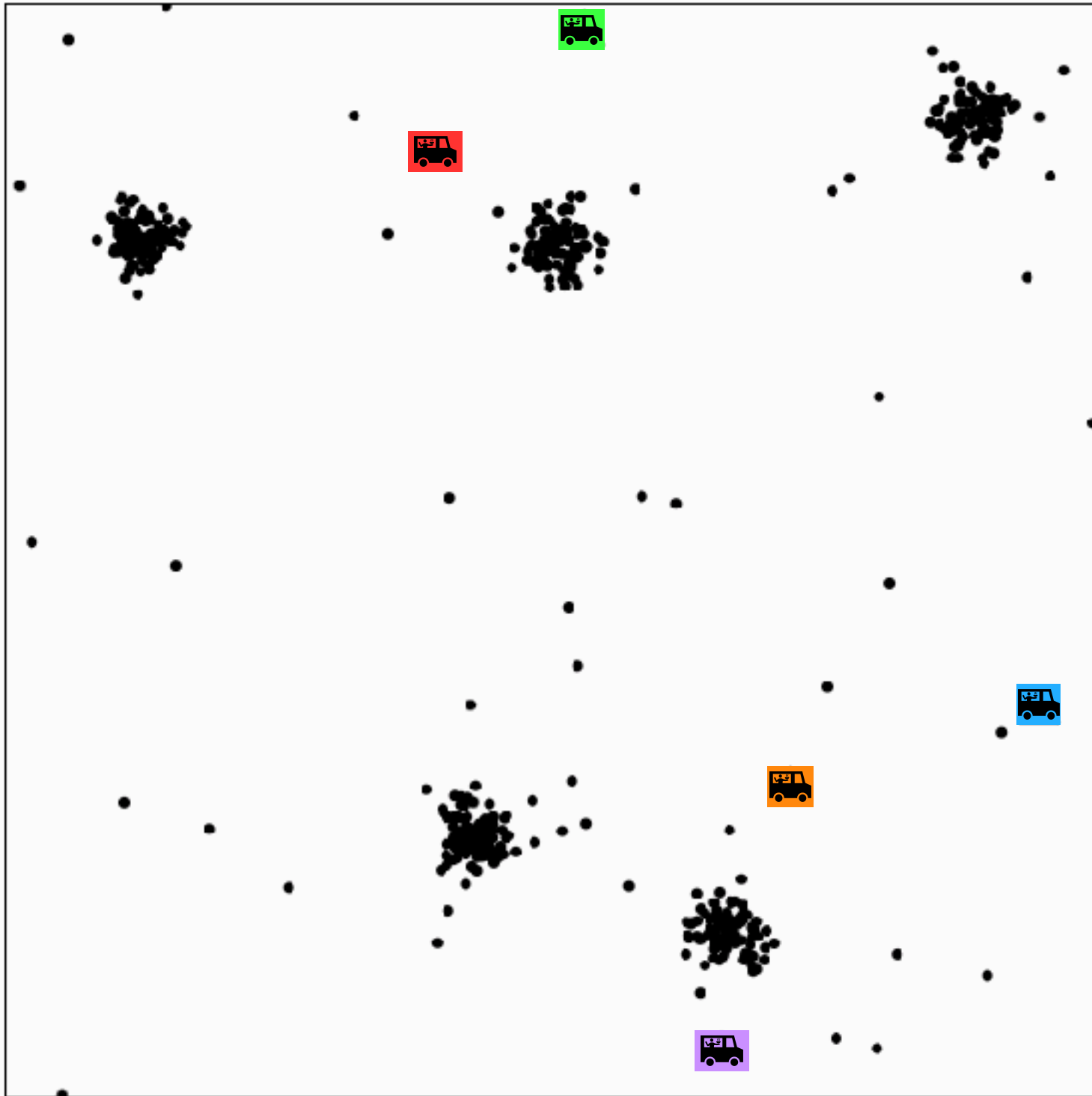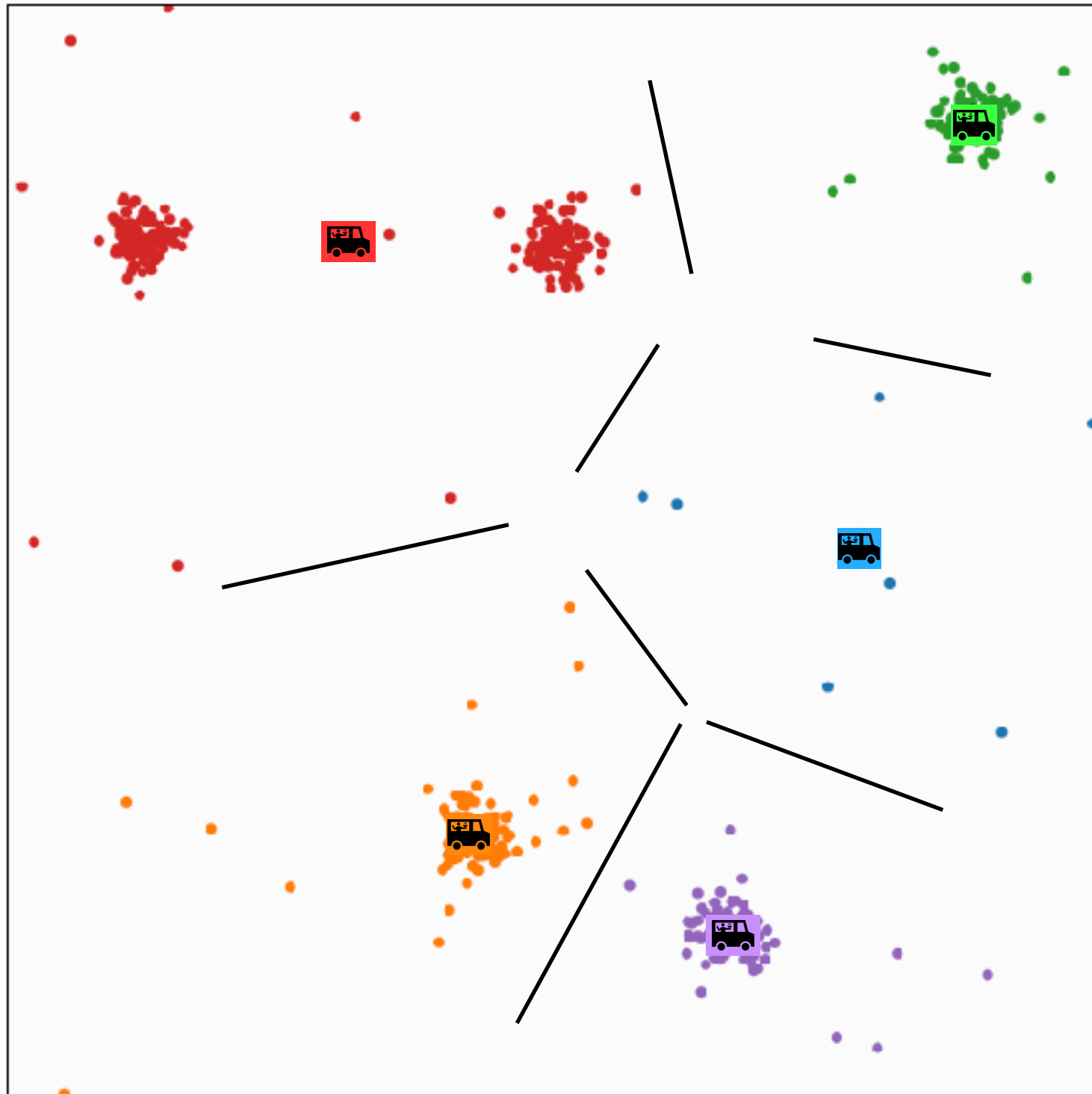# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!
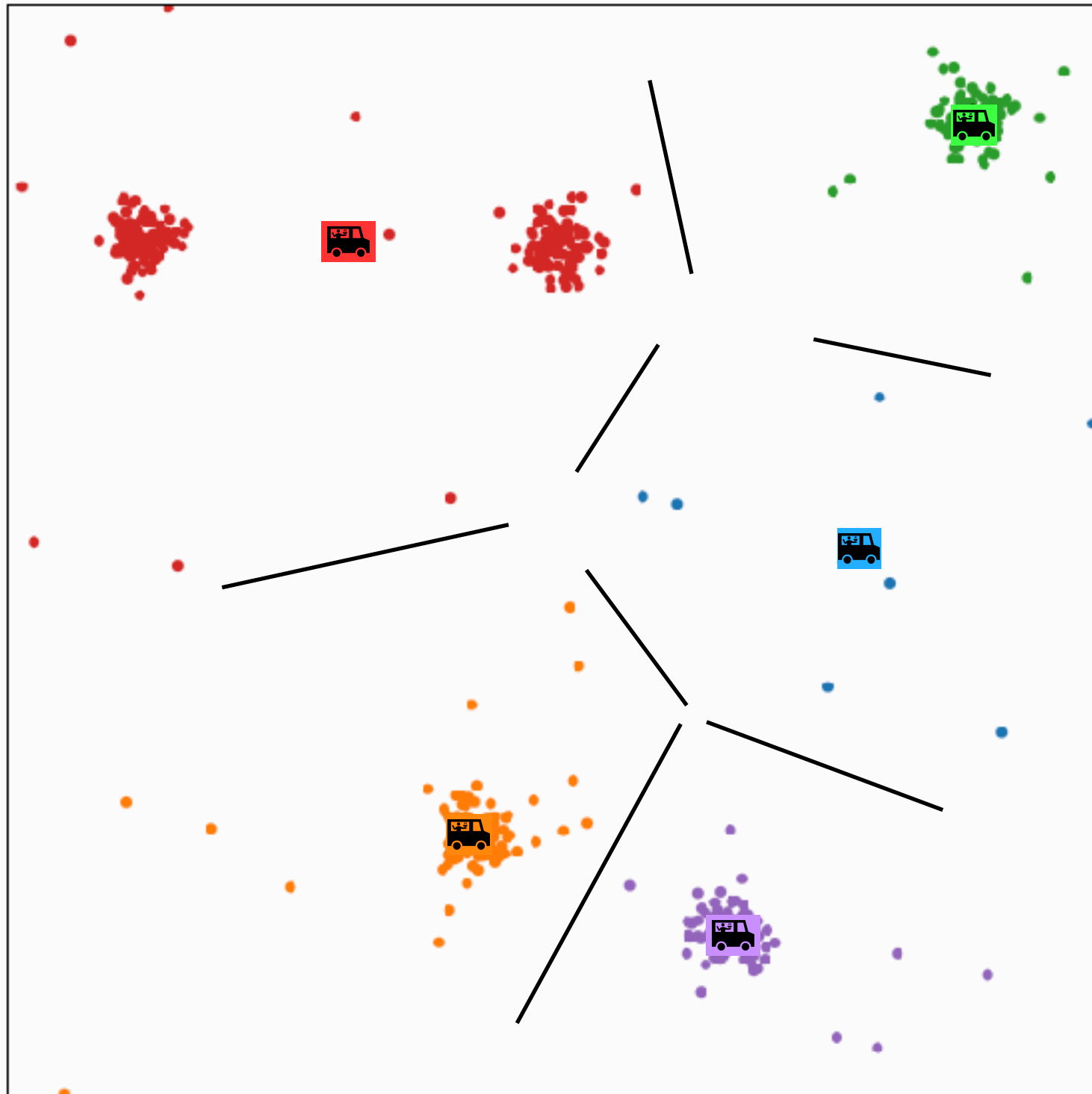
# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!

# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!

# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective
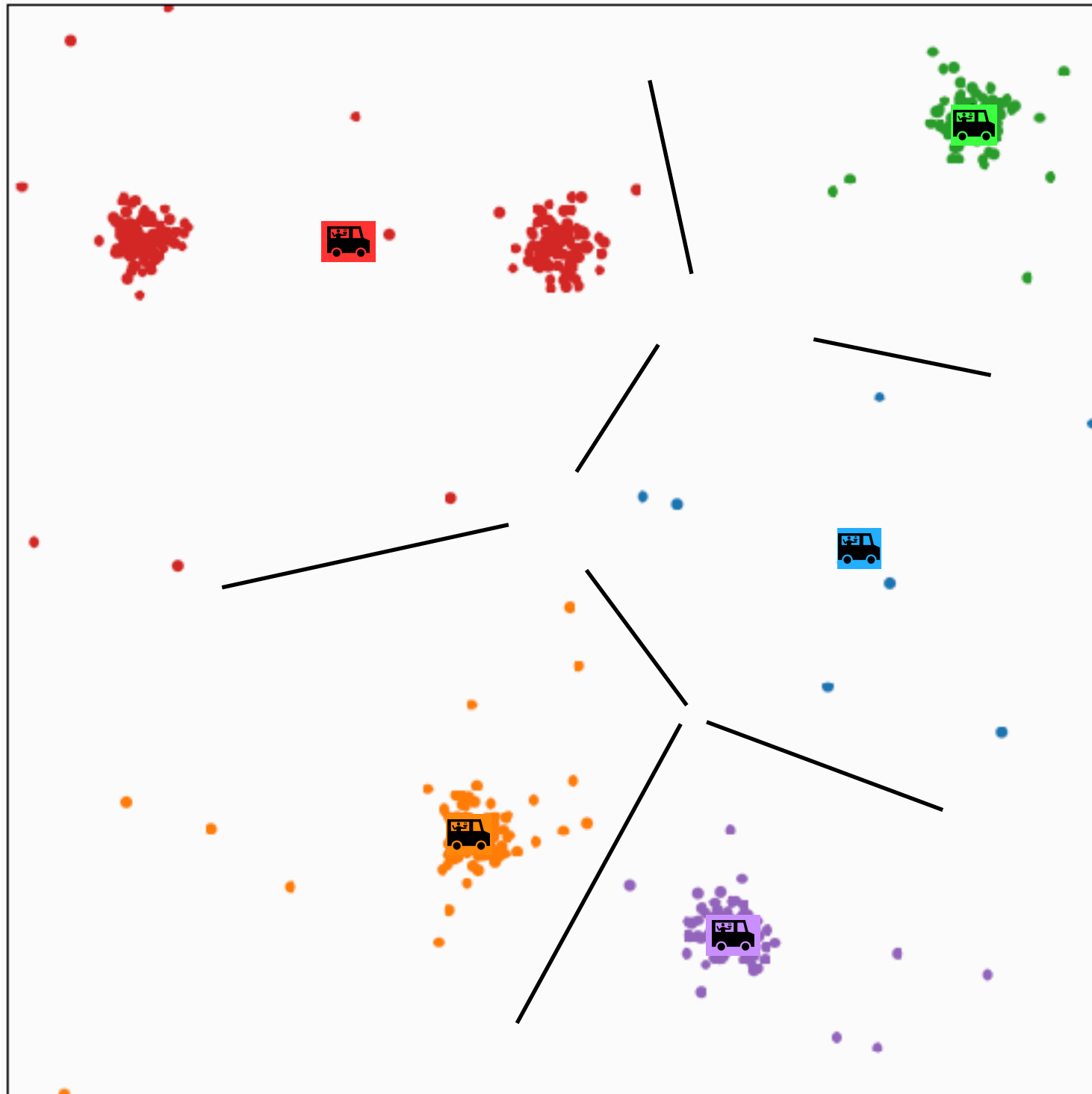- That local minimum could be bad!

# k-means algorithm: initialization
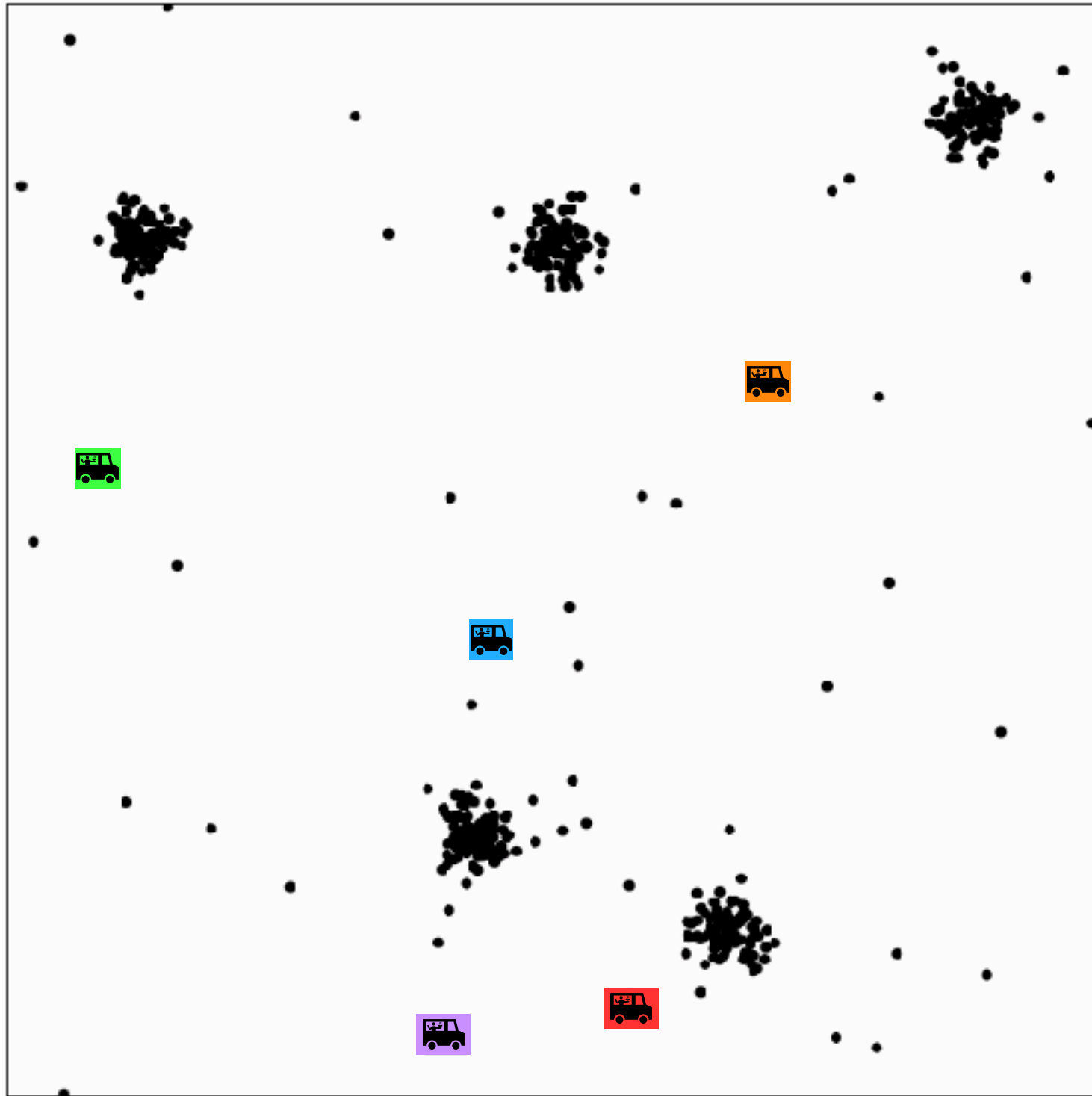


- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!

Is this clustering worse than the one we found before?

# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!

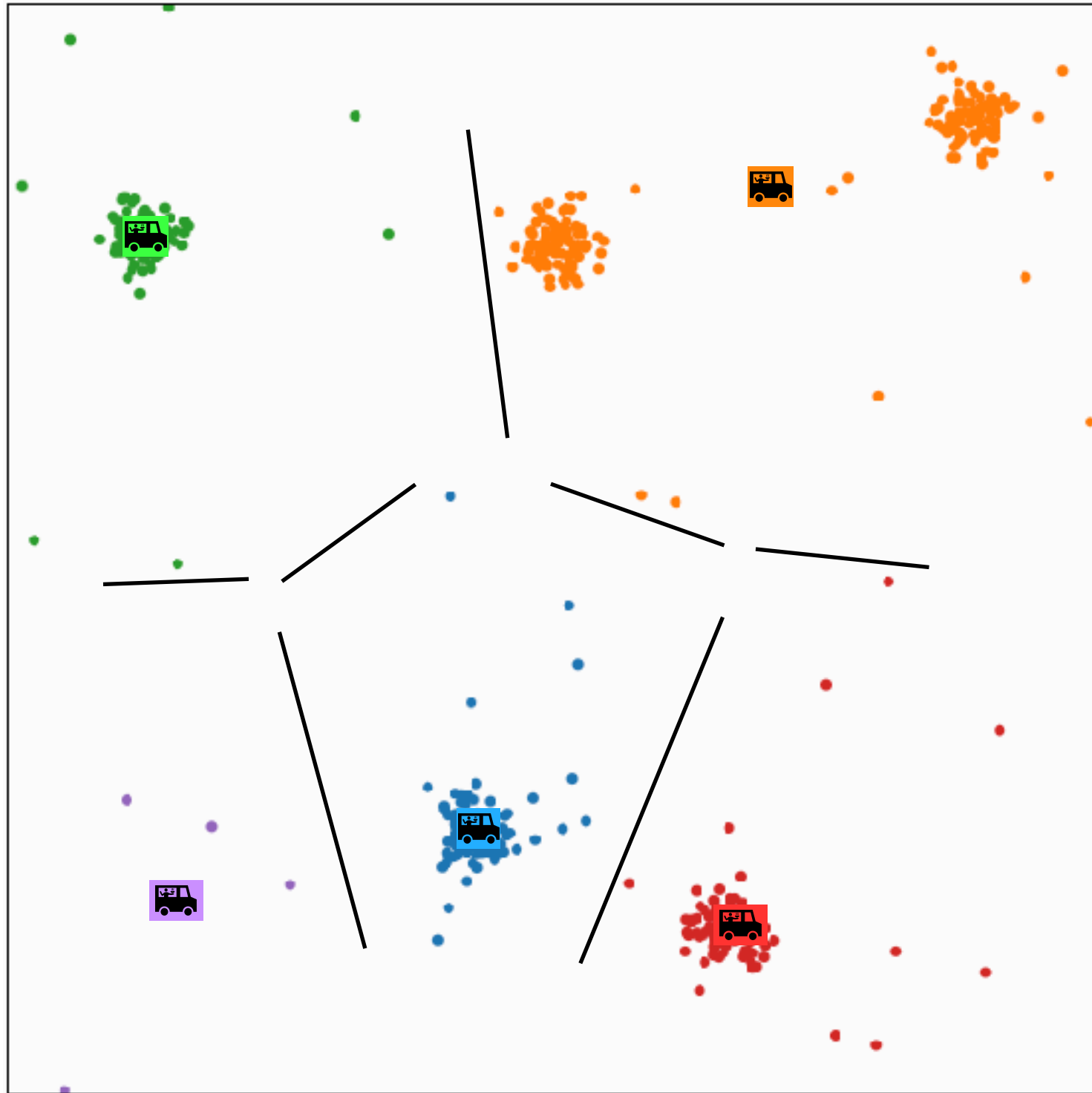Is this clustering worse than the one we found before?
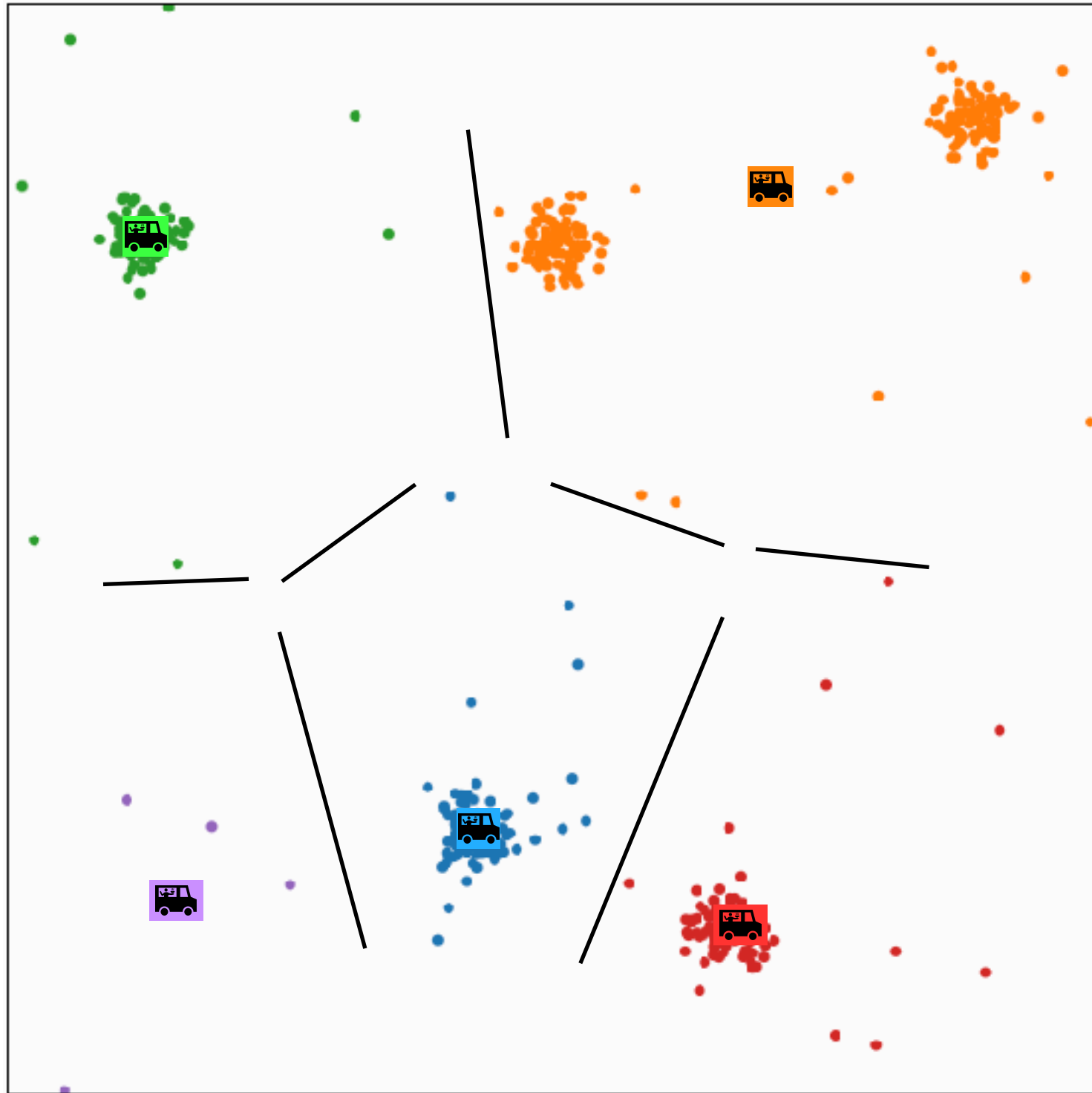
Why or why not?

# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!
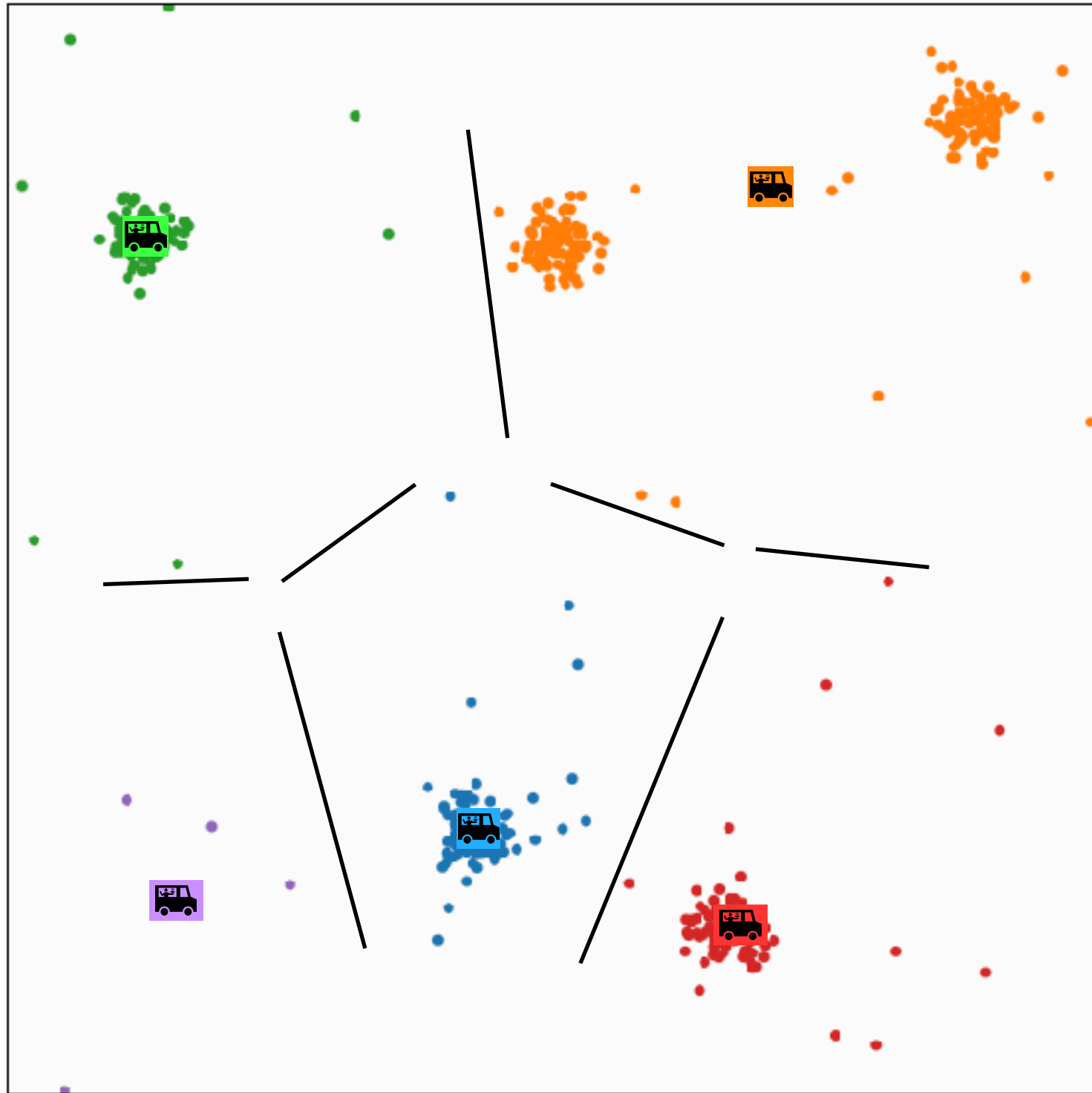
# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!

# k-means algorithm: initialization
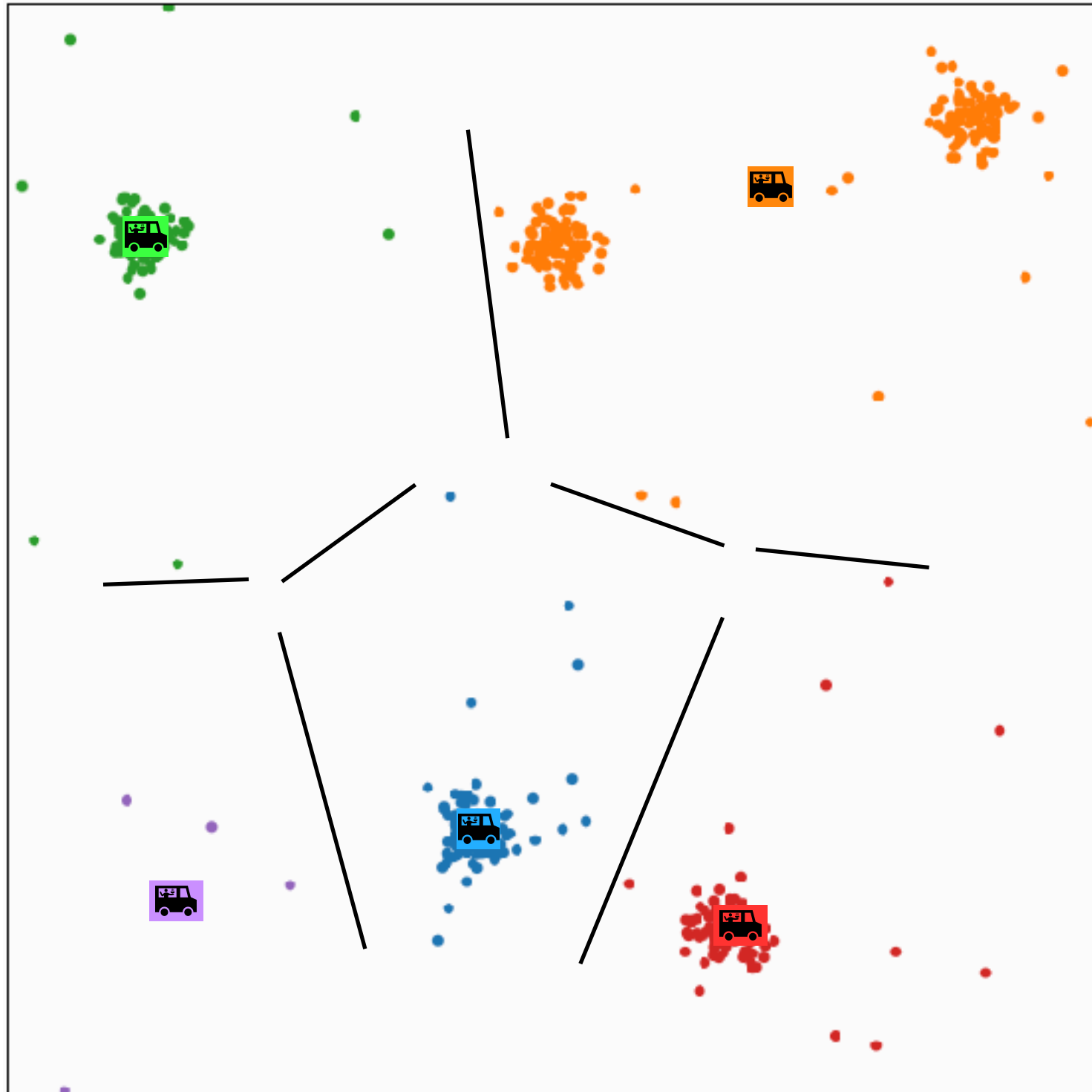


- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!

- The initialization can make a big difference
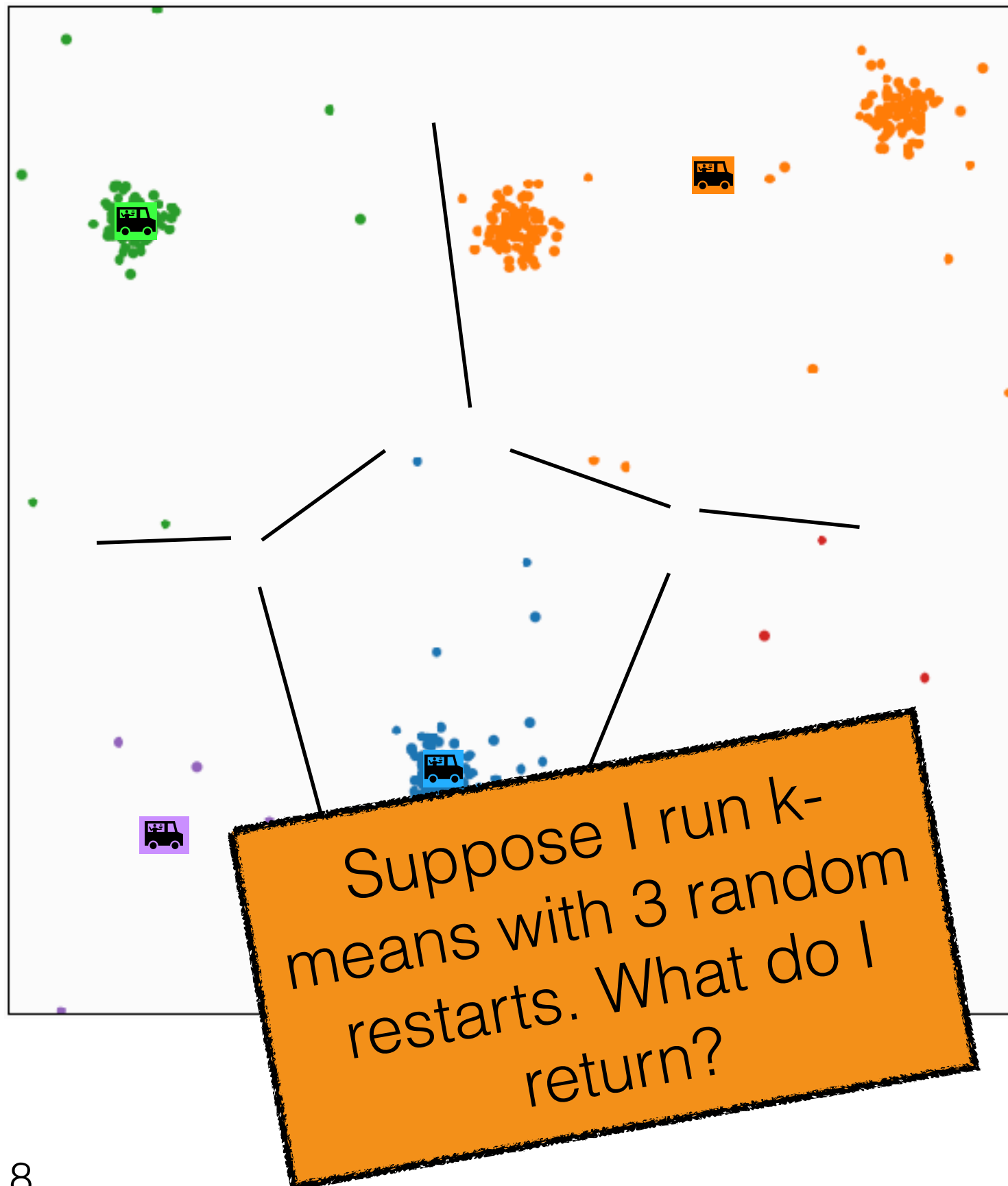
# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!

- The initialization can make a big difference

- Some options:
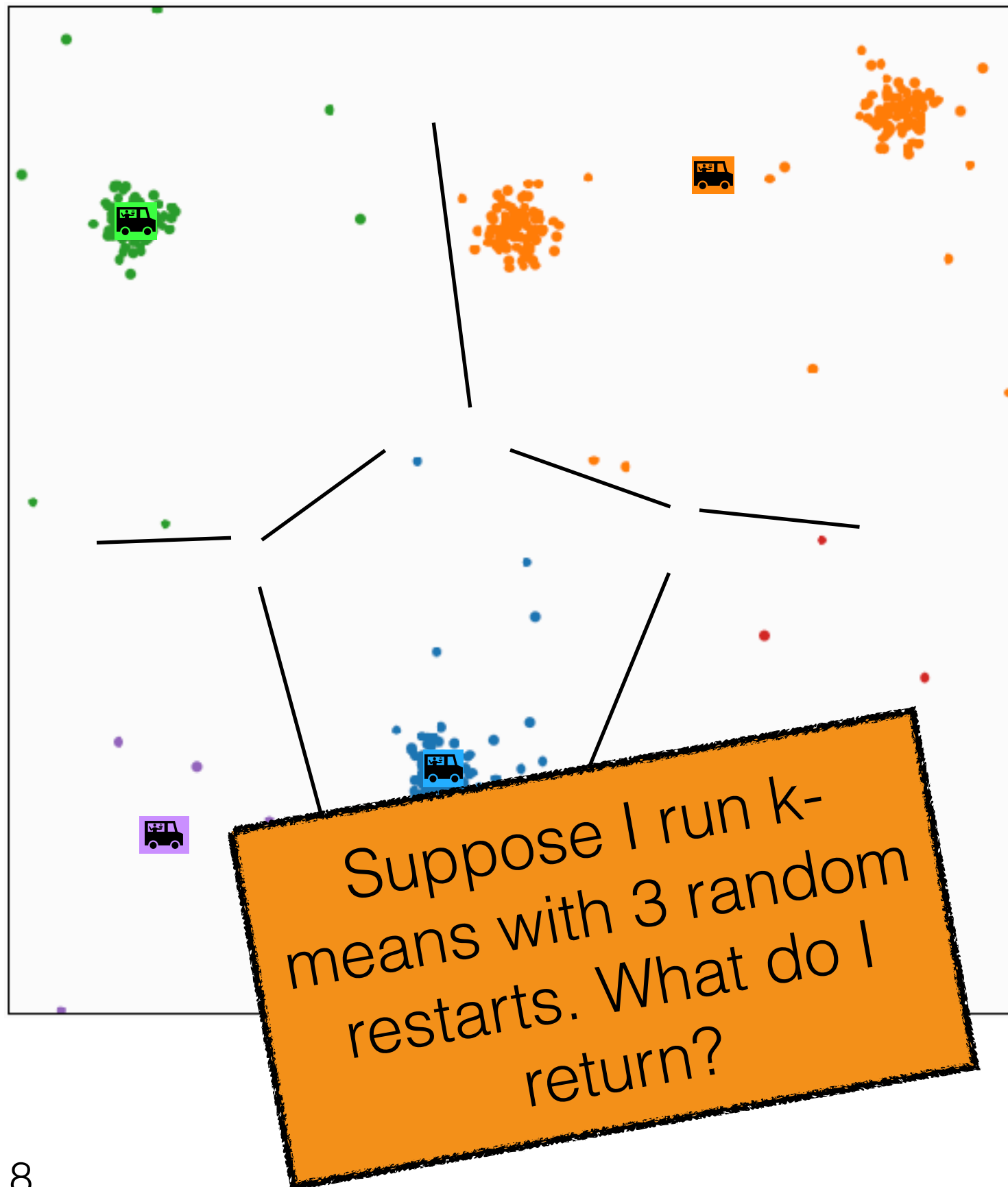
# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective
- That local minimum could be bad!
- The initialization can make a big difference
- Some options: random restarts

# k-means algorithm: initialization



- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!

- The initialization can make a big difference

- Some options: random restarts

Suppose I run k-means with 3 random restarts. What do I return?

# k-means algorithm: initialization



Suppose I run k-means with 3 random restarts. What do I return?

- **Theorem**. If run for enough outer iterations, the k-means algorithm will converge to a local minimum of the k-means objective

- That local minimum could be bad!

- The initialization can make a big difference
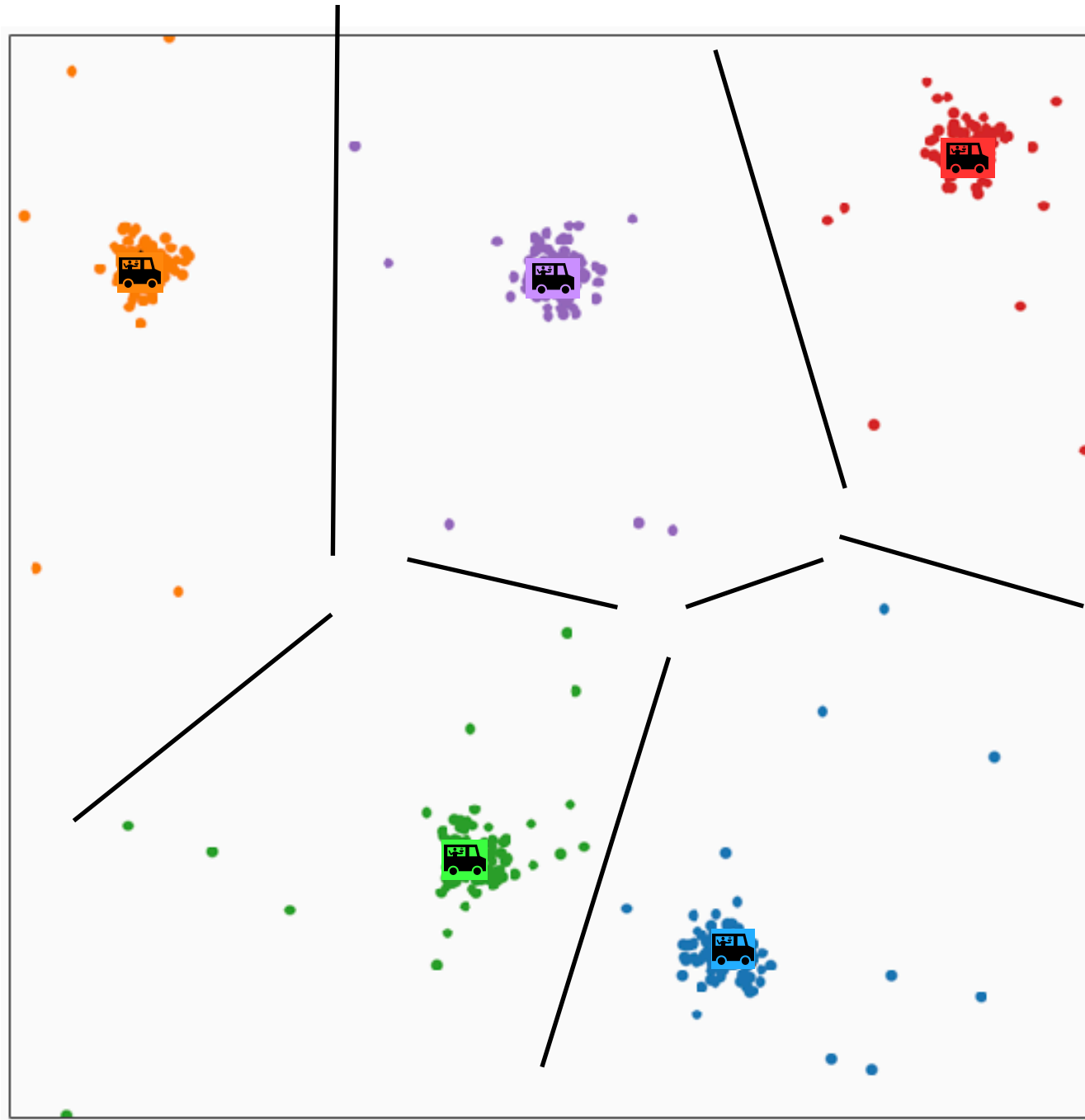
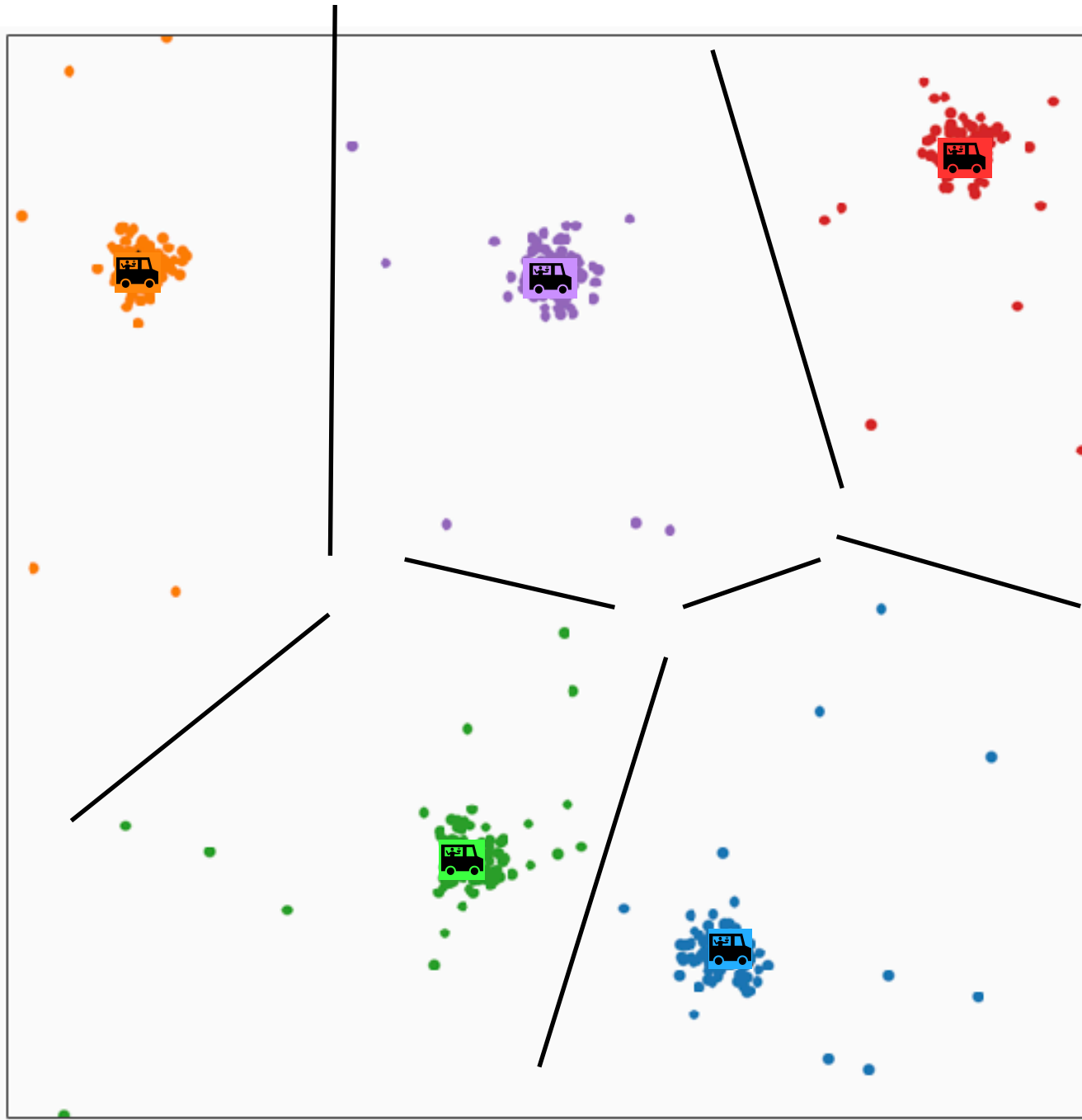- Some options: random restarts, k-means++

# k-means algorithm: effect of *k*
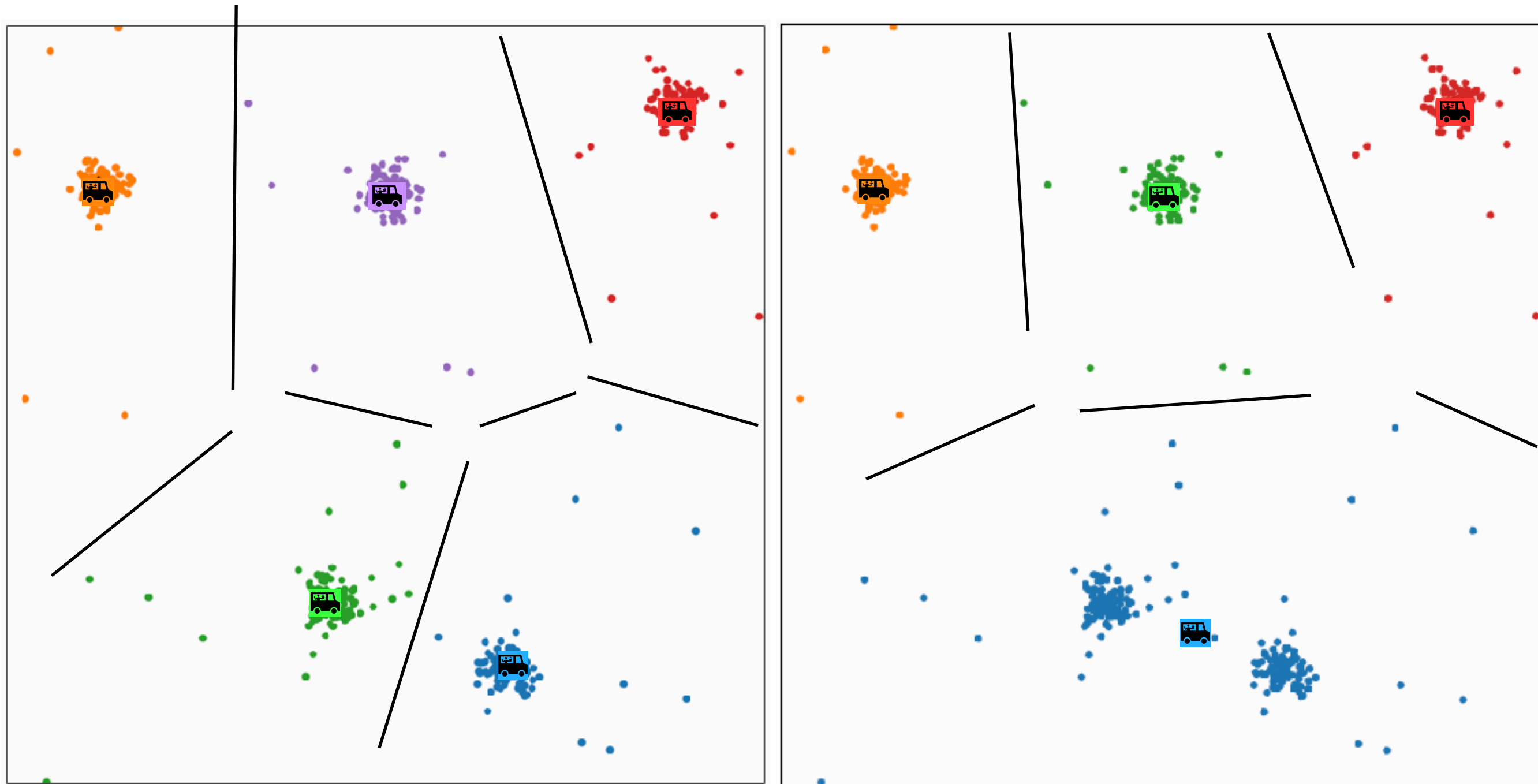
# k-means algorithm: effect of *k*

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]

# k-means algorithm: effect of *k*

- Different *k* will give us different results

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]

# k-means algorithm: effect of *k*

- Different *k* will give us different results

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]

# k-means algorithm: effect of *k*

- Different *k* will give us different results
- Larger *k* gets trucks closer to people

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]

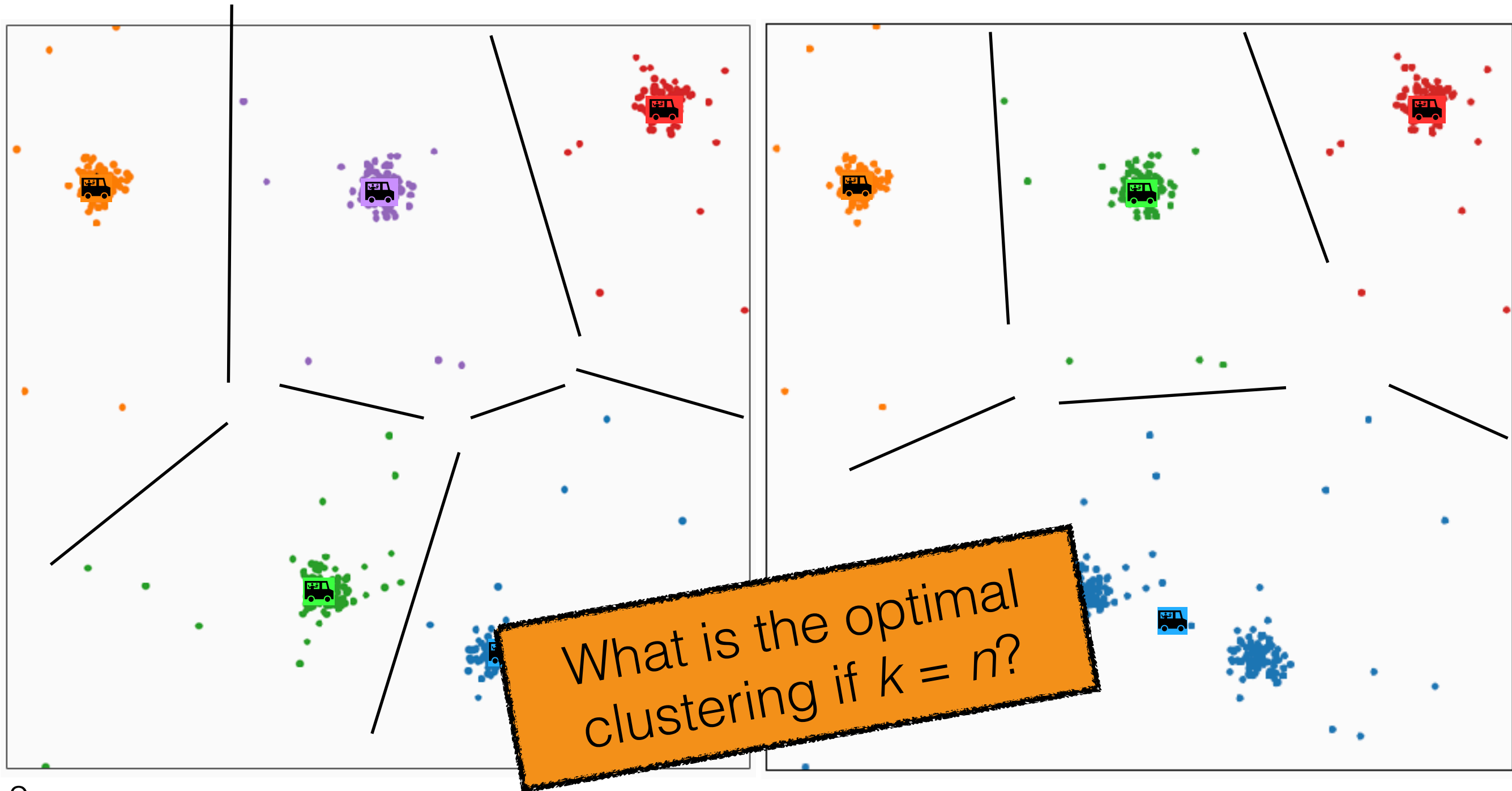# k-means algorithm: effect of *k*

- Different *k* will give us different results
- Larger *k* gets trucks closer to people



What is the optimal clustering if $k = n$?

[https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html]

# k-means algorithm: choosing *k*

# k-means algorithm: choosing *k*

- Sometimes we know *k*

# k-means algorithm: choosing *k*

- Sometimes we know *k*

# k-means algorithm: choosing *k*

- Sometimes we know *k*

# k-means algorithm: choosing *k*

- Sometimes we know *k*

# k-means algorithm: choosing *k*

- Sometimes we know *k*

# k-means algorithm: choosing *k*

- Sometimes we know *k*

# k-means algorithm: choosing *k*

- Sometimes we know *k*



- Sometimes we'd like to choose/learn *k*

# k-means algorithm: choosing *k*

- Sometimes we know *k*



- Sometimes we'd like to choose/learn *k*
  - Can't just minimize the k-means objective over *k* too

# k-means algorithm: choosing *k*

- Sometimes we know *k*

- Sometimes we'd like to choose/learn *k*
  - Can't just minimize the k-means objective over *k* too

$$\sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

# k-means algorithm: choosing *k*

- Sometimes we know *k*



- Sometimes we'd like to choose/learn *k*
  - Can't just minimize the k-means objective over *k* too

$$\arg\min_{y,\mu} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

# k-means algorithm: choosing *k*

- Sometimes we know *k*



- Sometimes we'd like to choose/learn *k*
  - Can't just minimize the k-means objective over *k* too

$$\arg\min_{y,\mu,k} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

# k-means algorithm: choosing *k*

- Sometimes we know *k*



- Sometimes we'd like to choose/learn *k*

  - Can't just minimize the k-means objective over *k* too

$$\arg\min_{y,\mu,k} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

Why not?

# k-means algorithm: choosing *k*

- Sometimes we know *k*



- Sometimes we'd like to choose/learn *k*
  - Can't just minimize the k-means objective over *k* too

$$\arg\min_{y,\mu,\boxed{k}} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

# k-means algorithm: choosing *k*

- Sometimes we know *k*

- Sometimes we'd like to choose/learn *k*
  - Can't just minimize the k-means objective over *k* too

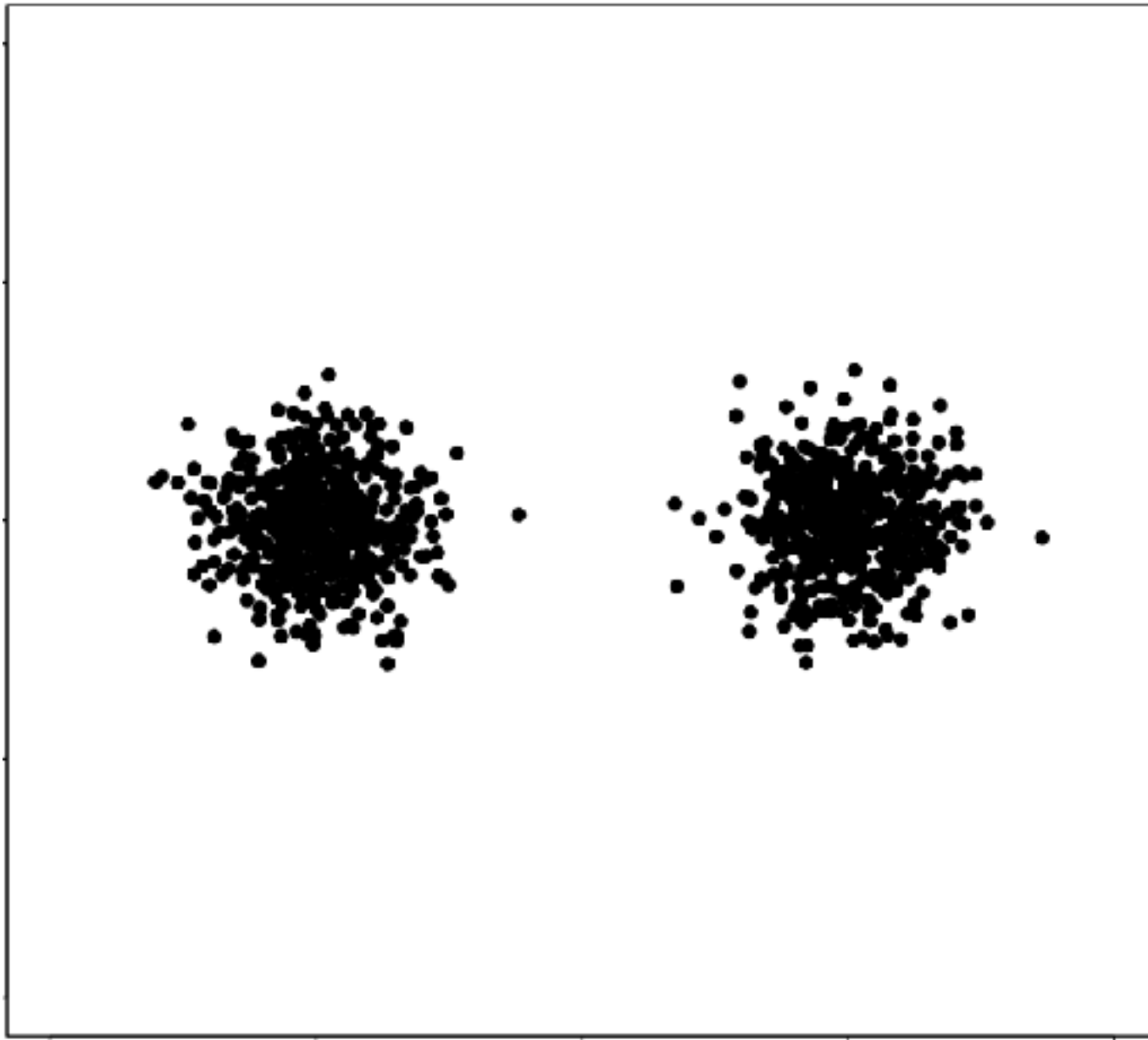$$\arg\min_{y,\mu,k} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

- How to choose *k* depends on what you'd like to do

# k-means algorithm: choosing *k*

- Sometimes we know *k*

- Sometimes we'd like to choose/learn *k*
  - Can't just minimize the k-means objective over *k* too

$$\arg\min_{y,\mu,k} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2$$

- How to choose *k* depends on what you'd like to do
  - E.g. cost-benefit trade-off

# k-means algorithm: choosing *k*

- Sometimes we know *k*



- Sometimes we'd like to choose/learn *k*
  - Can't just minimize the k-means objective over *k* too

$$\arg\min_{y,\mu,k} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2 + \mathrm{cost}(k)$$

- How to choose *k* depends on what you'd like to do
  - E.g. cost-benefit trade-off

# k-means algorithm: choosing *k*

- Sometimes we know *k*



- Sometimes we'd like to choose/learn *k*
  - Can't just minimize the k-means objective over *k* too

$$\arg\min_{y,\mu,k} \sum_{j=1}^{k} \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}\|x^{(i)} - \mu^{(j)}\|_2^2 + \mathrm{cost}(k)$$

- How to choose *k* depends on what you'd like to do
  - E.g. cost-benefit trade-off
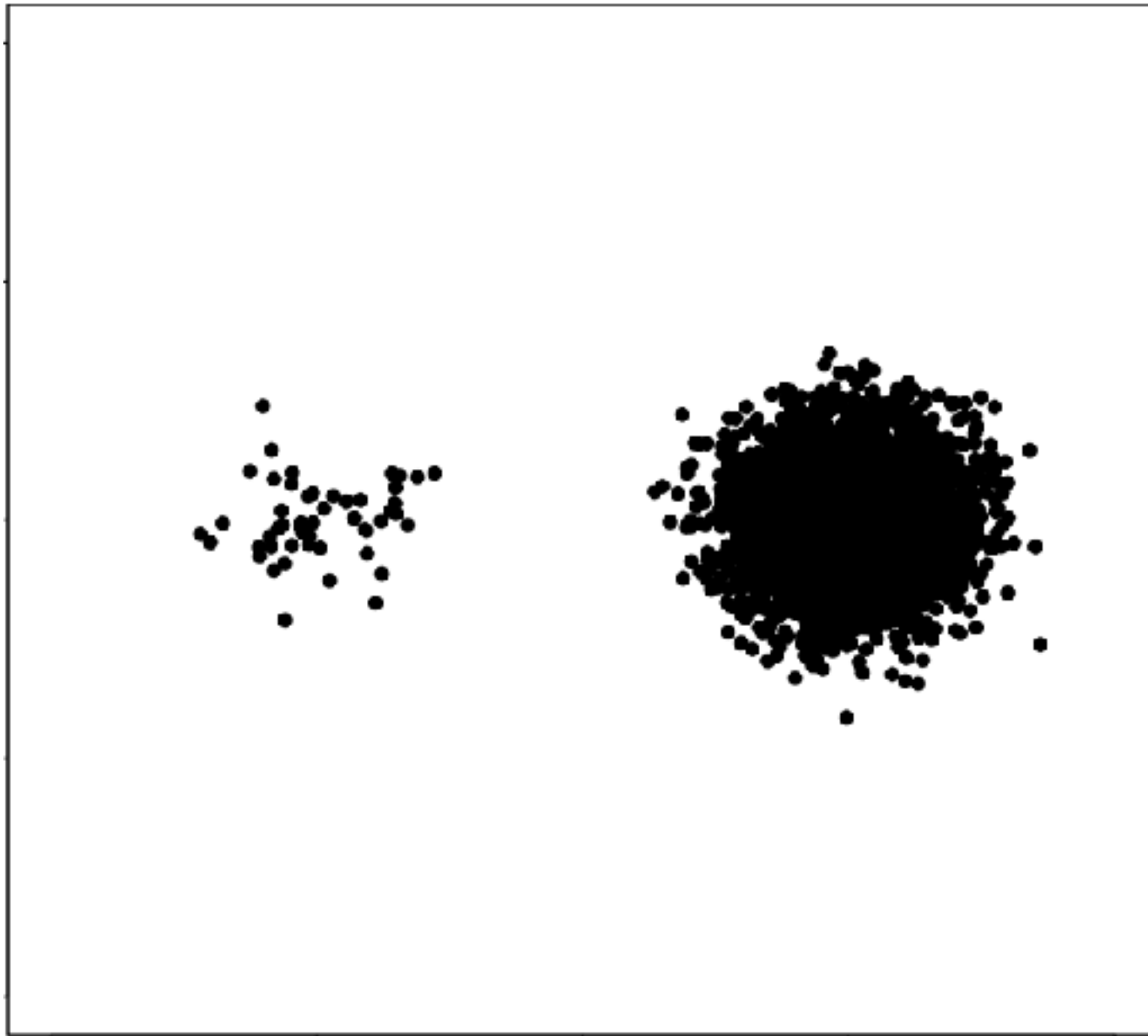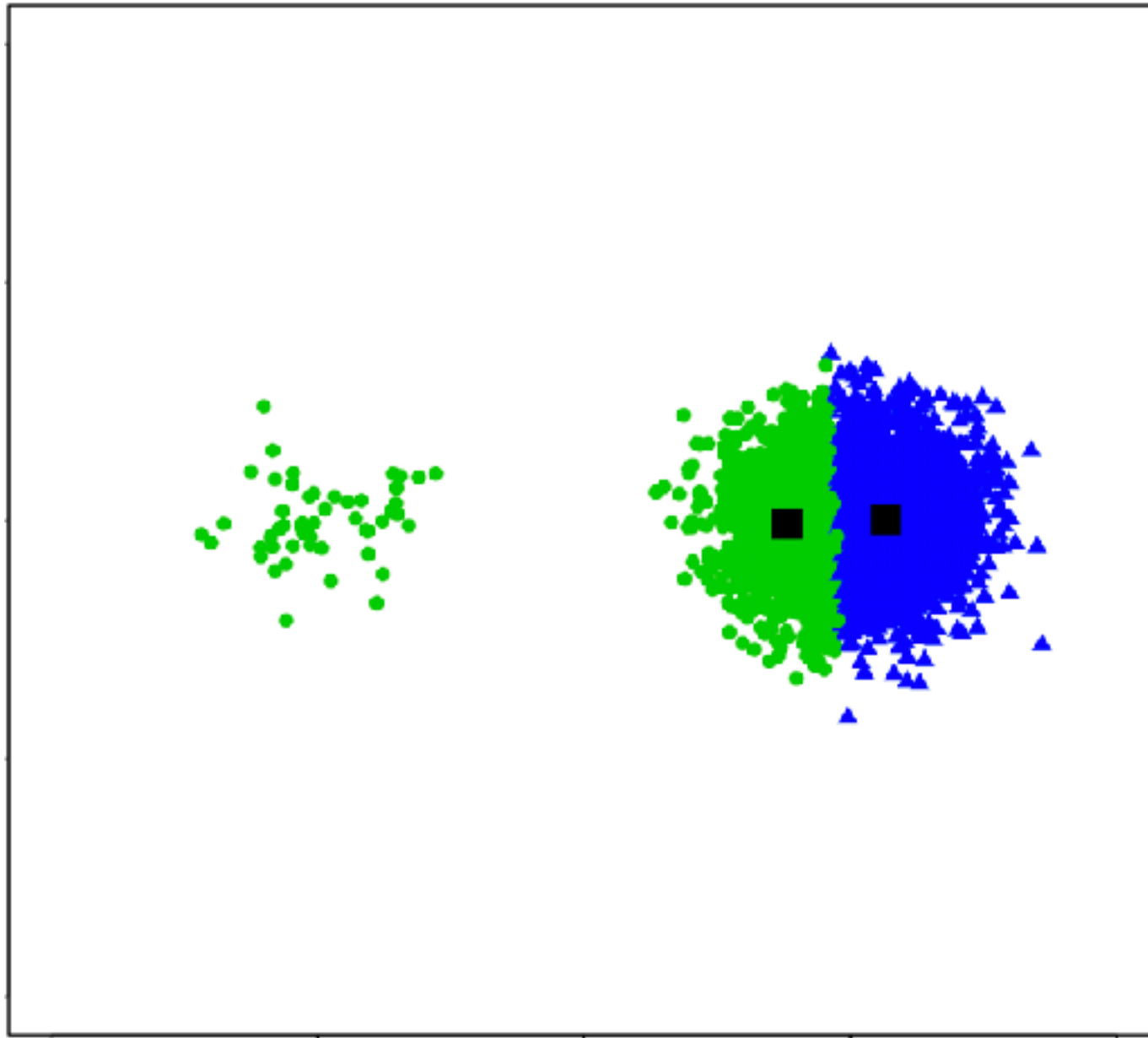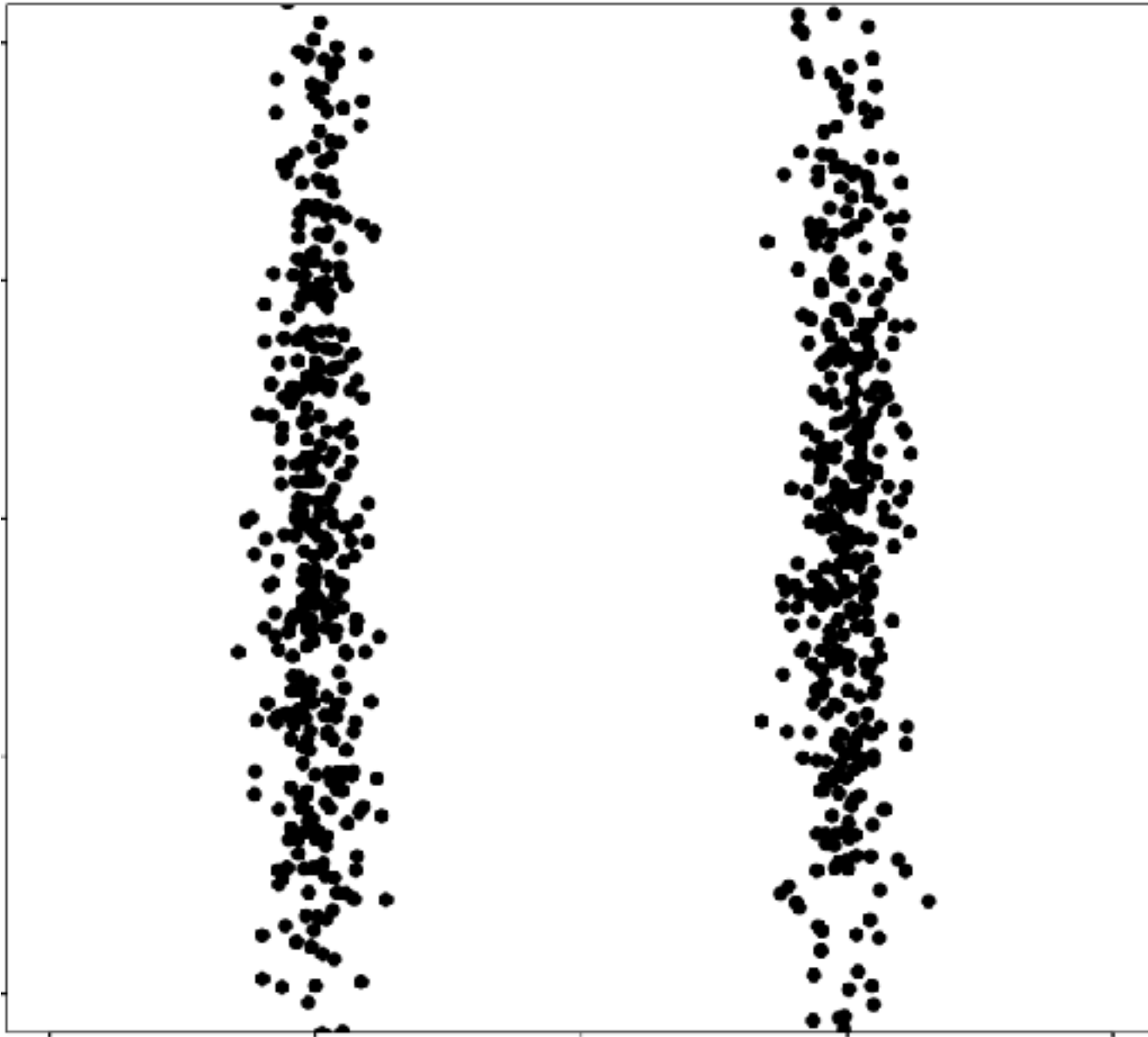  - Often no single "right answer"

# Cluster shape

# Cluster shape

- k-means works well for well-separated circular clusters of the same size

# Cluster shape



- k-means works well for well-separated circular clusters of the same size

# Cluster shape



- k-means works well for well-separated circular clusters of the same size
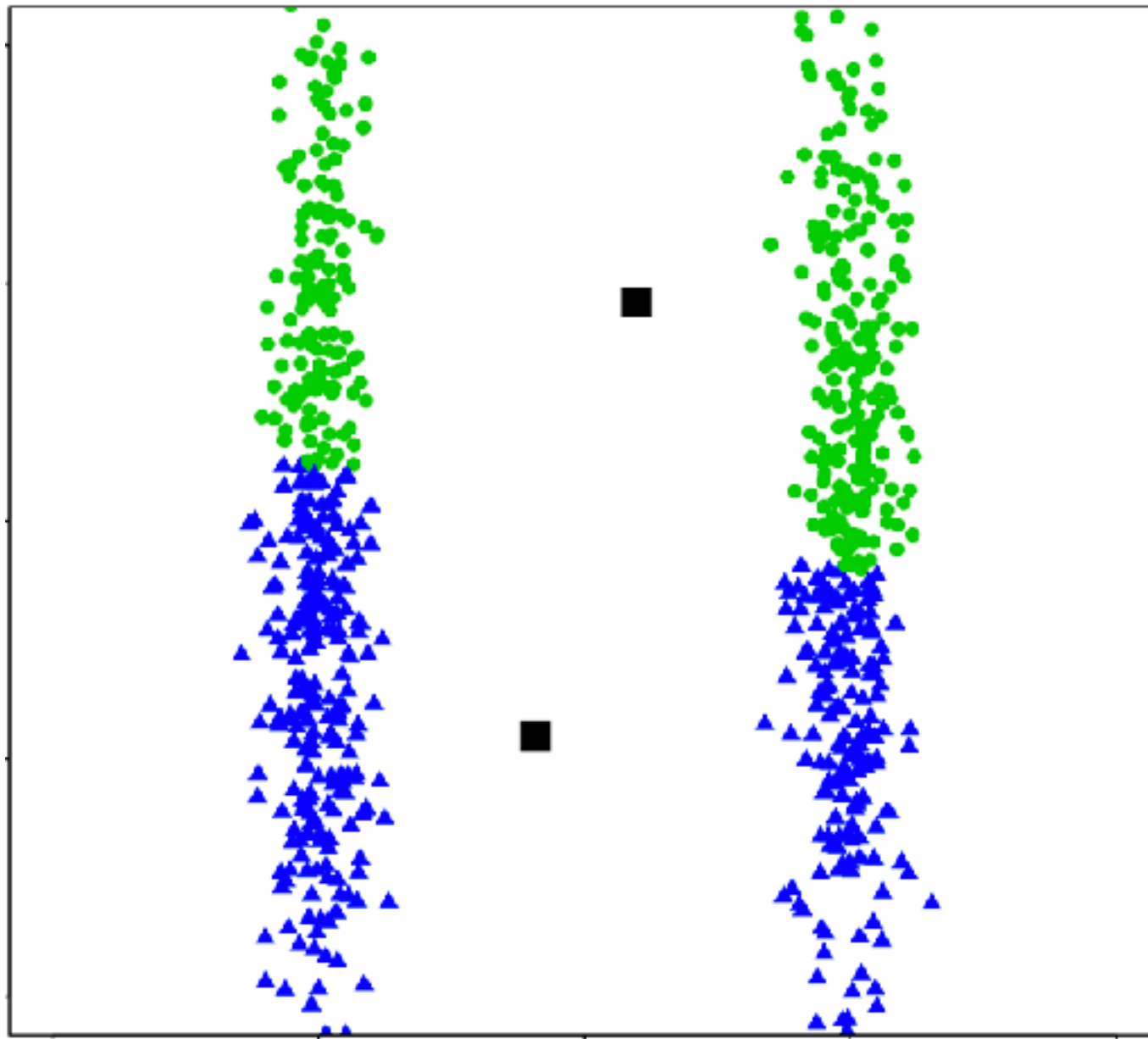
# Cluster shape



- k-means works well for well-separated circular clusters of the **same size**
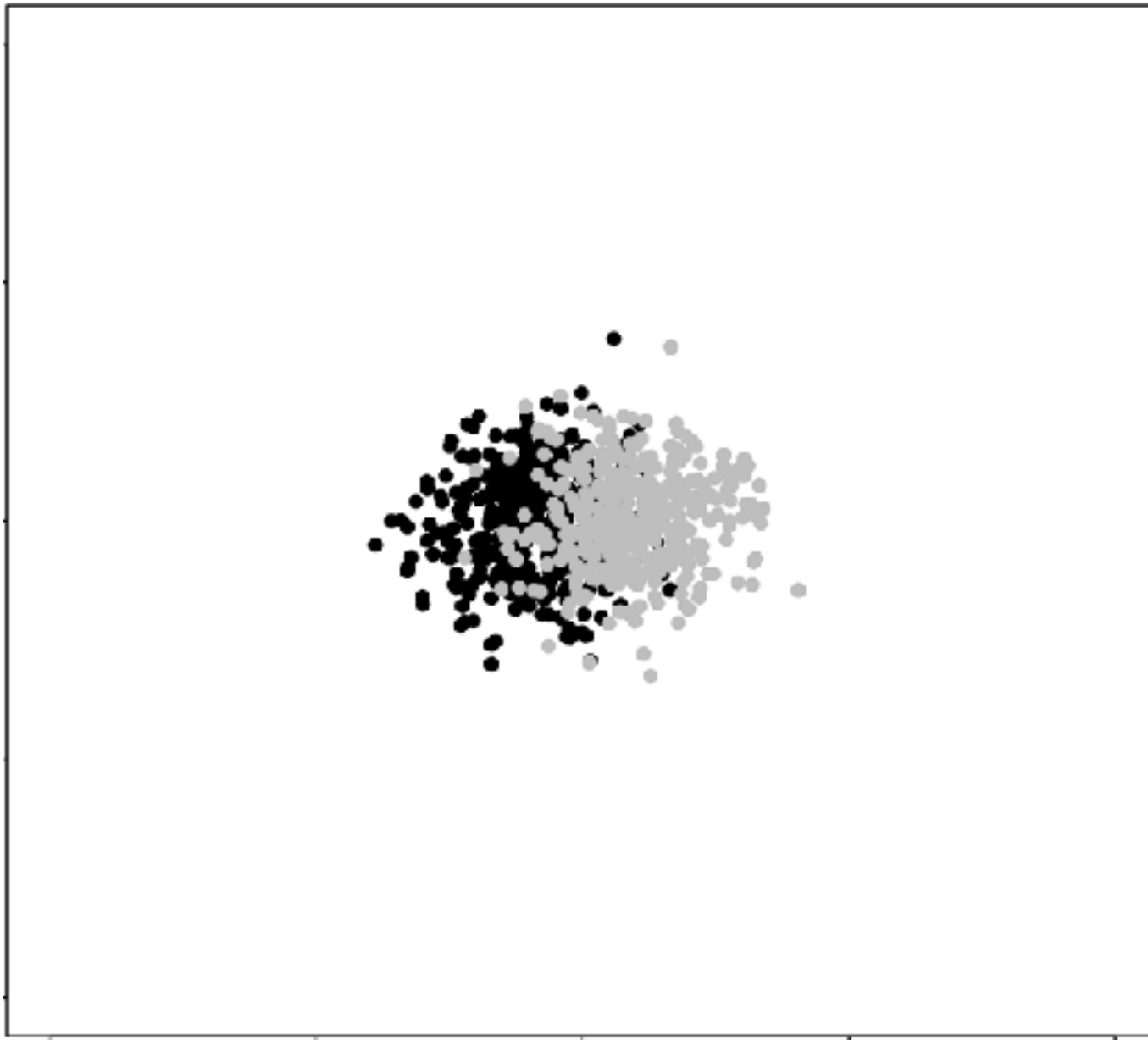
# Cluster shape



- k-means works well for well-separated circular clusters of the **same size**

# Cluster shape



- k-means works well for well-separated circular clusters of the **same size**

# Cluster shape



- k-means works well for well-separated circular clusters of the **same size**

# Cluster shape



- k-means works well for well-separated circular clusters of the **same size**

# Cluster shape



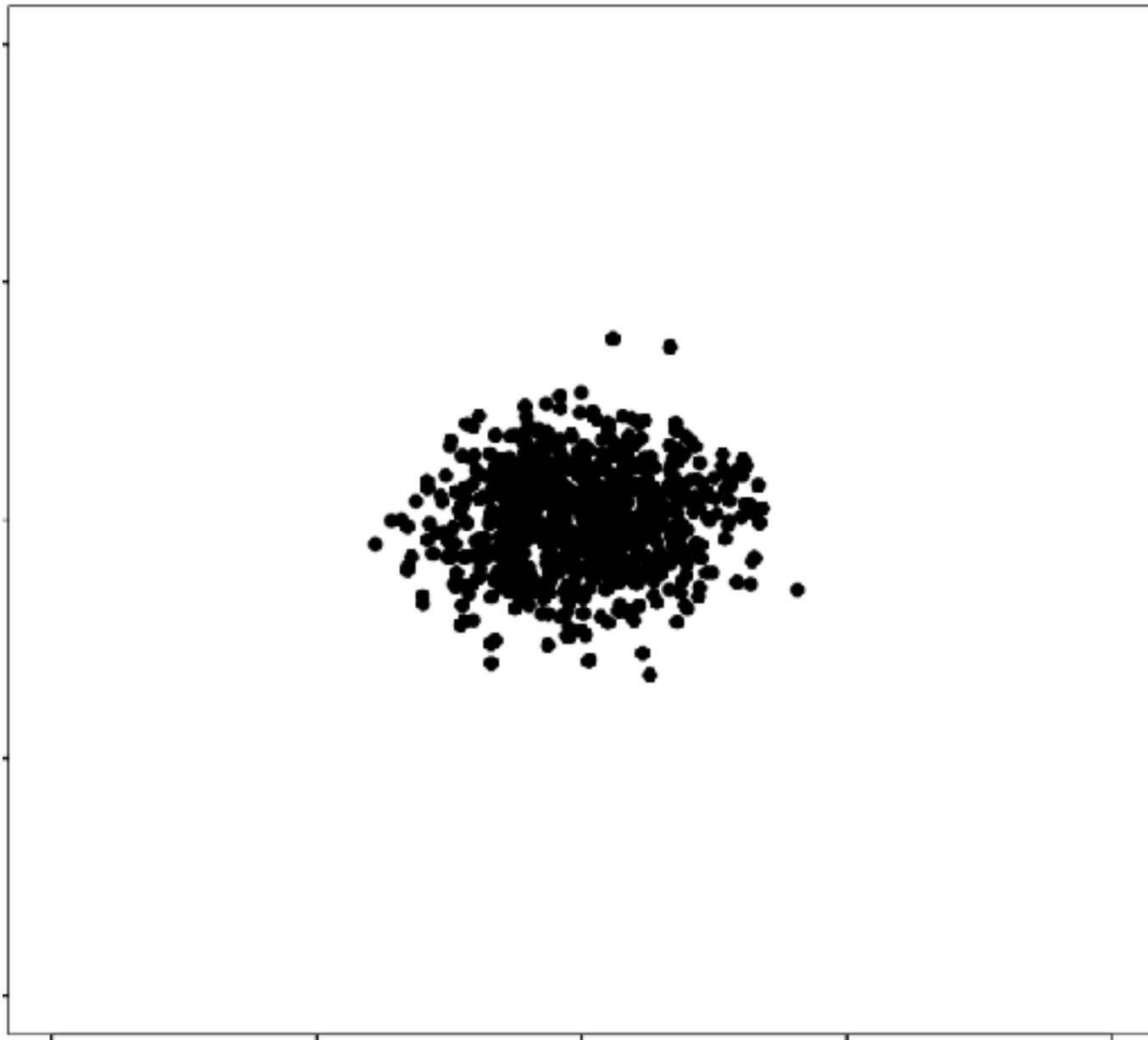- k-means works well for well-separated **circular** clusters of the same size

# Cluster shape



- k-means works well for well-separated **circular** clusters of the same size
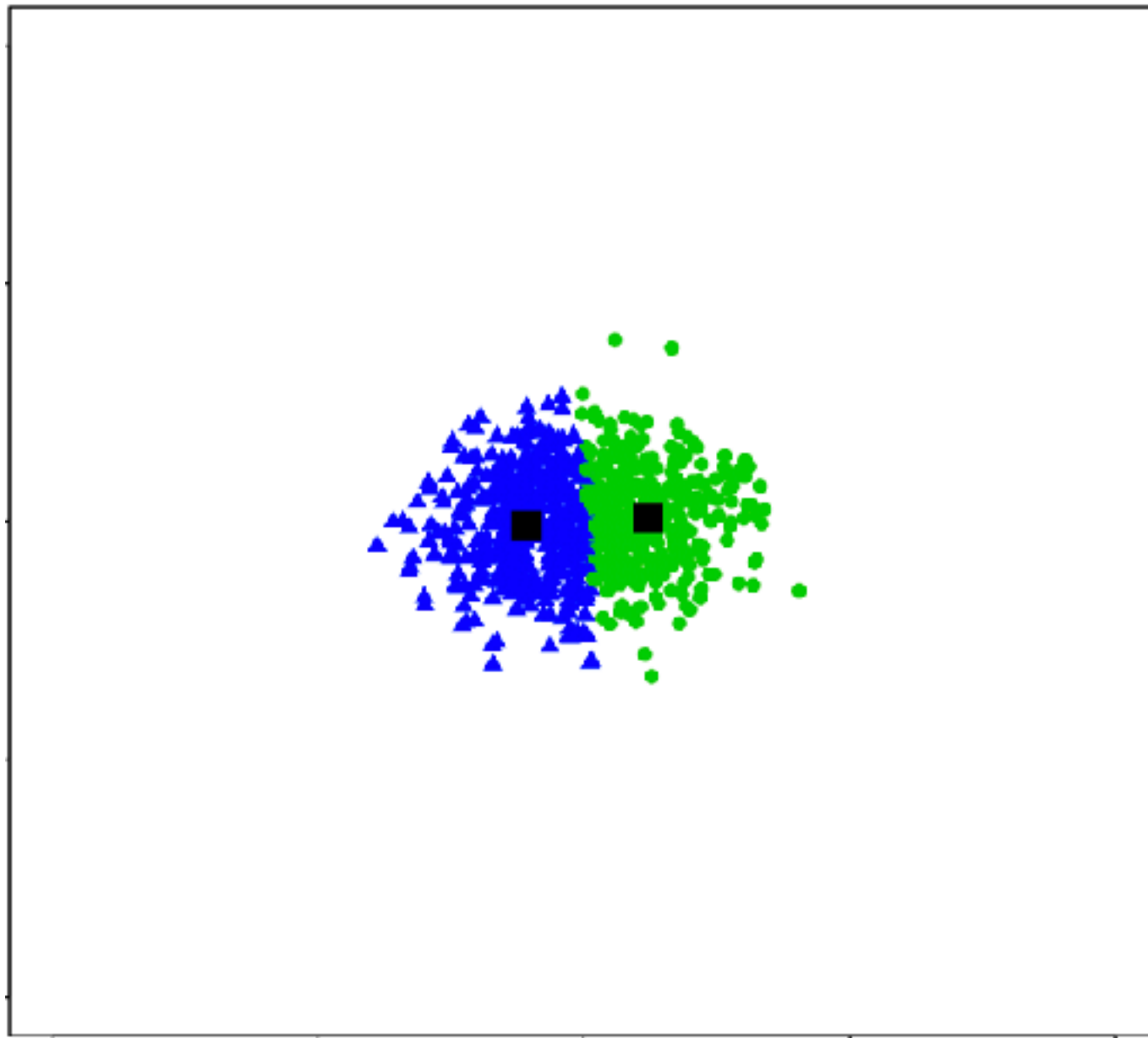
# Cluster shape



- k-means works well for **well-separated** circular clusters of the same size

# Cluster shape



- k-means works well for **well-separated** circular clusters of the same size
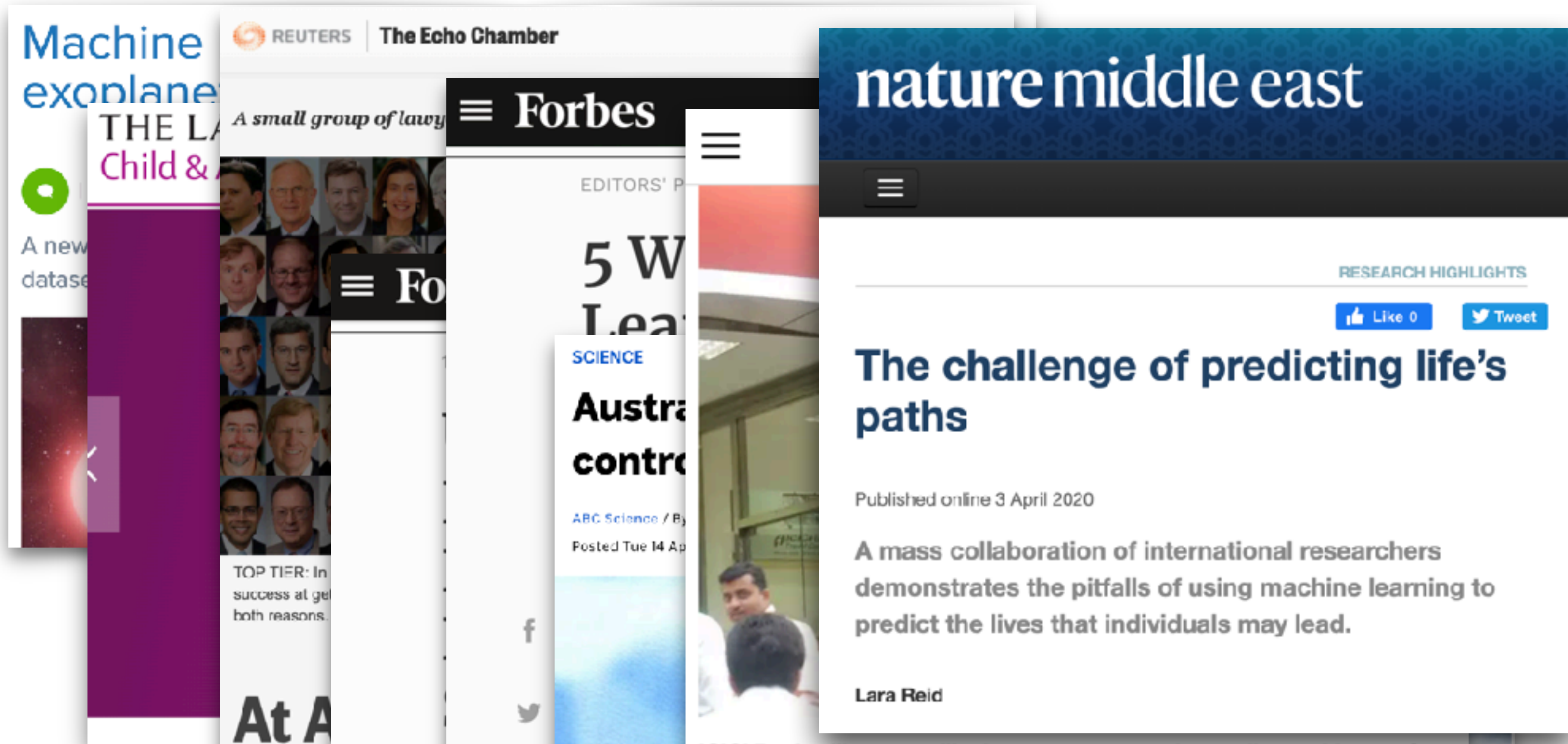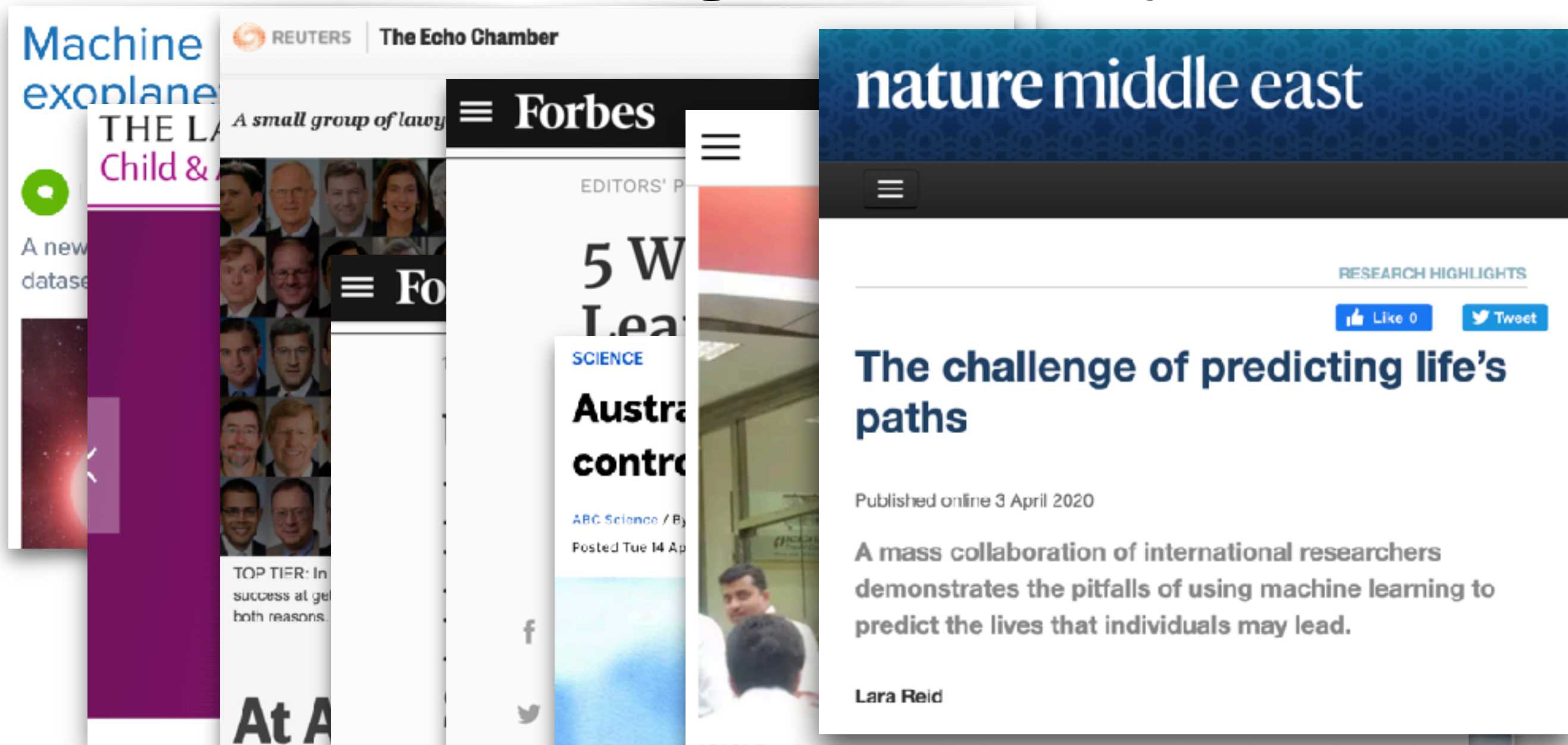
# Cluster shape



- k-means works well for **well-separated** circular clusters of the same size

# Machine learning (ML): why & what

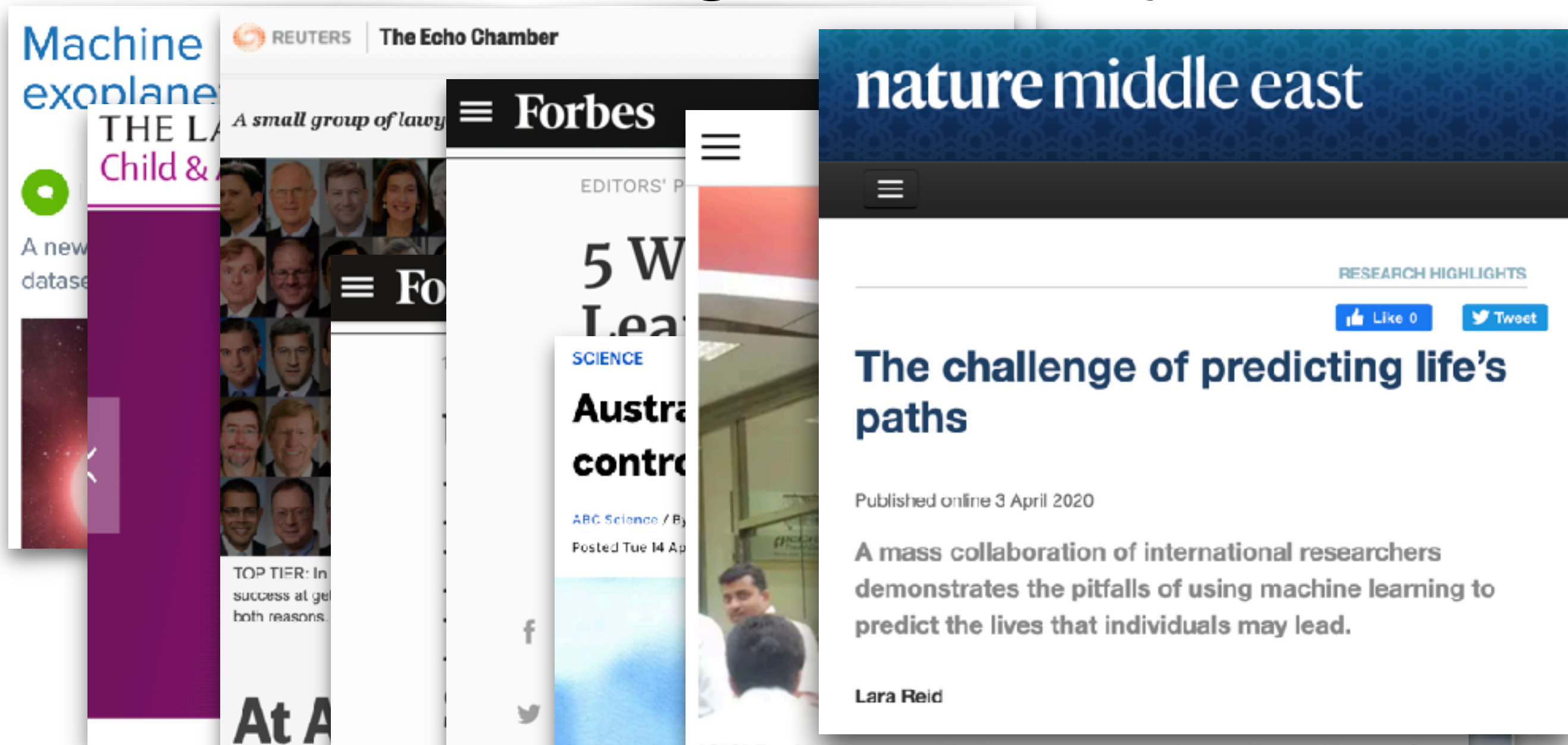# Machine learning (ML): why & what
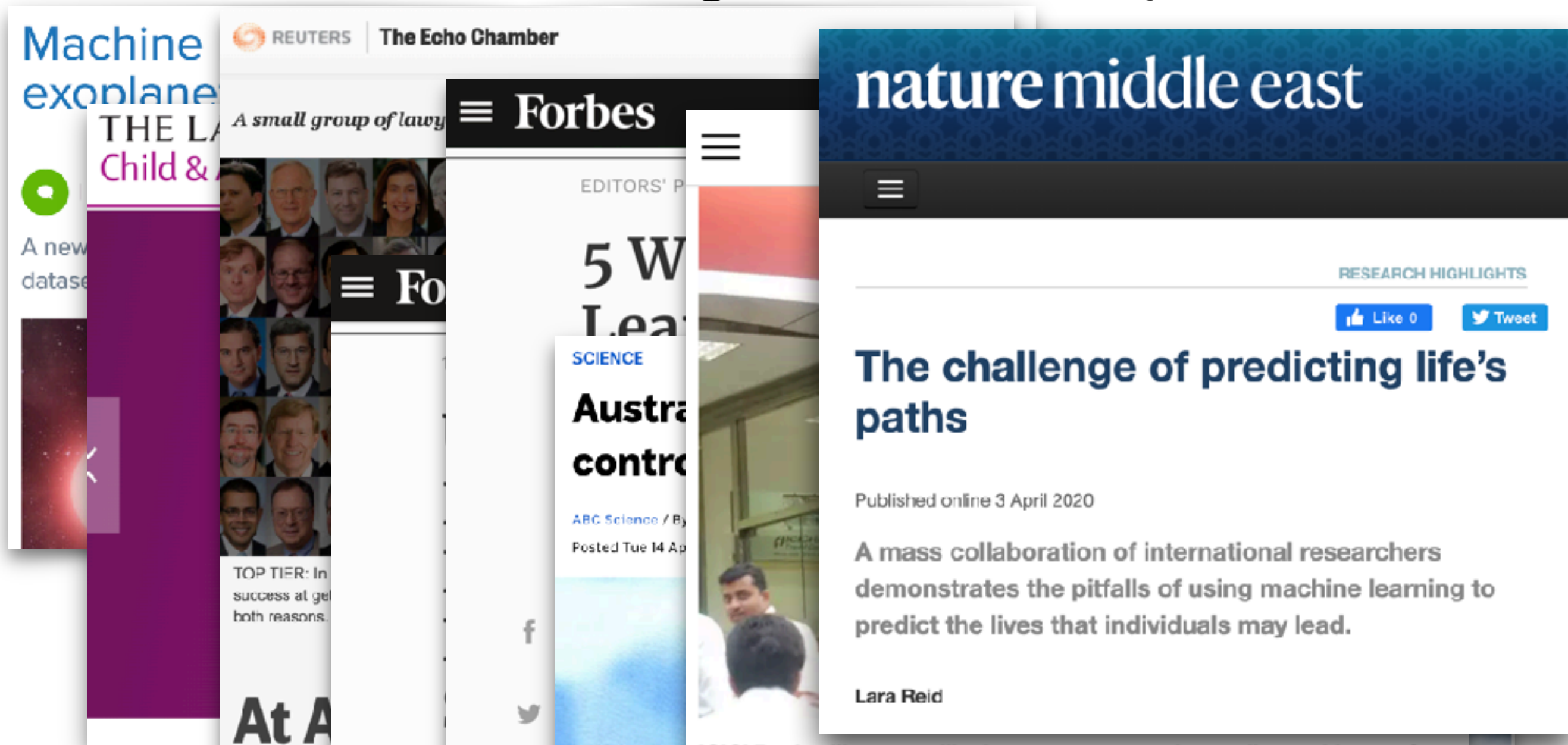
# Machine learning (ML): why & what



- **What is ML?** A set of methods for making decisions from data. (See the rest of the course!)

# Machine learning (ML): why & what



- **What is ML?** A set of methods for making decisions from data. (See the rest of the course!)
- **Why study ML?** To apply; to understand; to evaluate

# Machine learning (ML): why & what



- **What is ML?** A set of methods for making decisions from data. (See the rest of the course!)
- **Why study ML?** To apply; to understand; to evaluate
- **Notes:** ML is not magic. ML is built on math.