



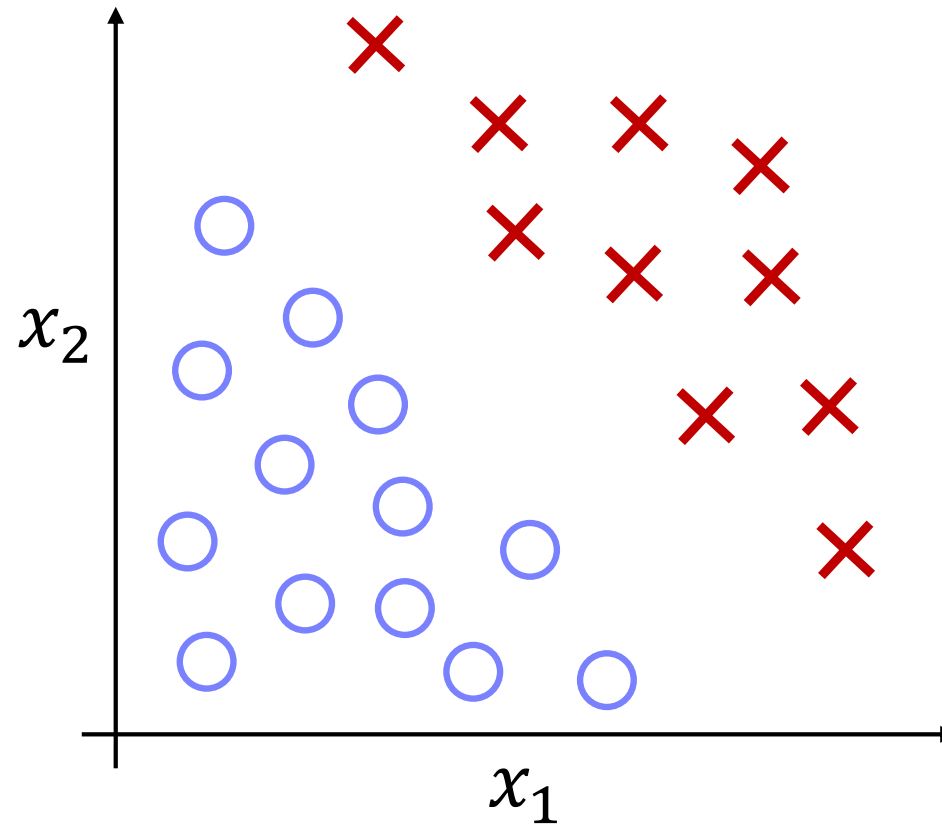
Machine Learning 11

Kihyun Shin
DMSE, HBNU

Clustering

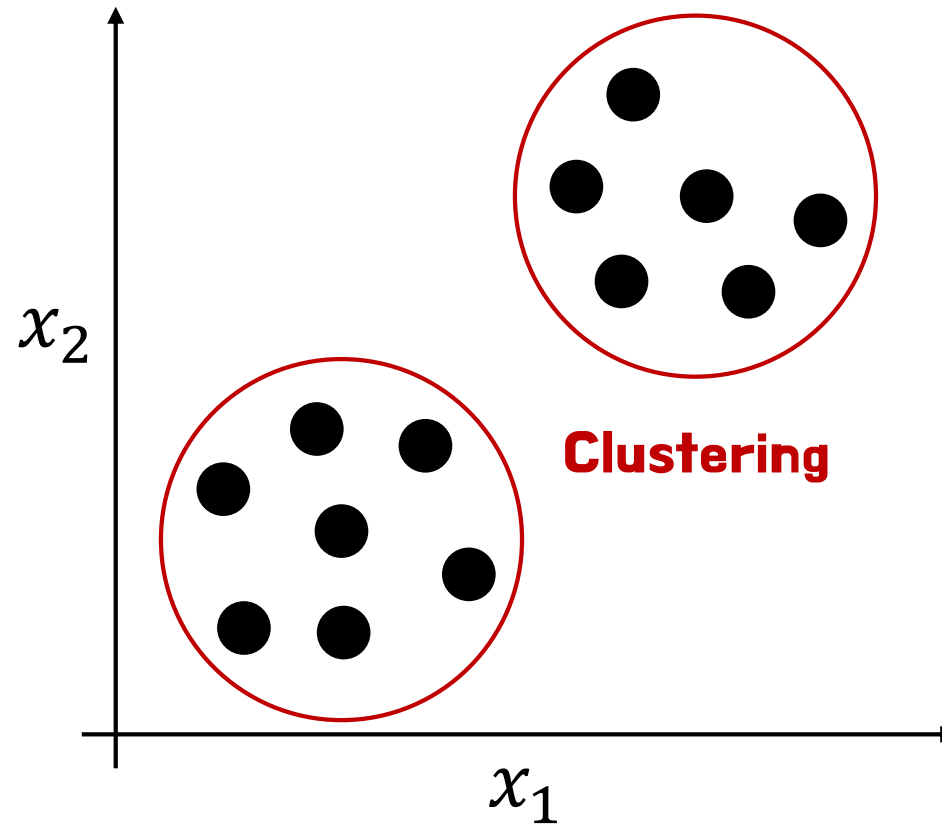


Supervised learning



Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$

Unsupervised learning

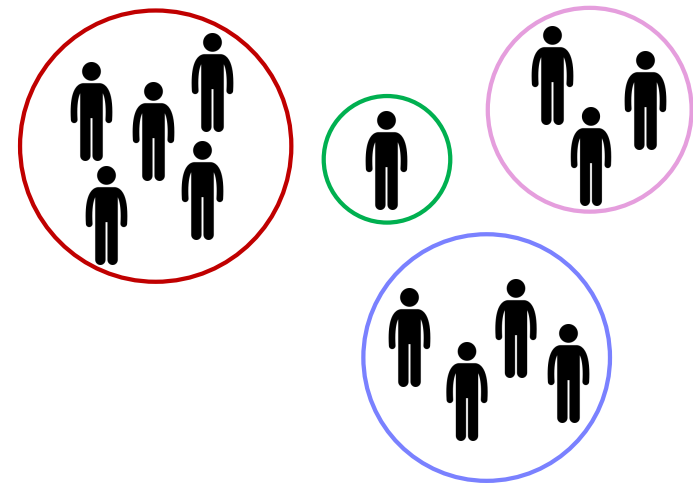


Training set: $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$

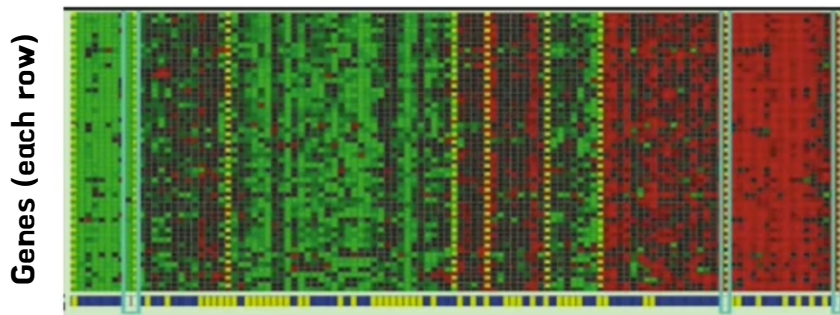
Applications of clustering



Grouping similar news



Market segmentation



Individuals (each column)

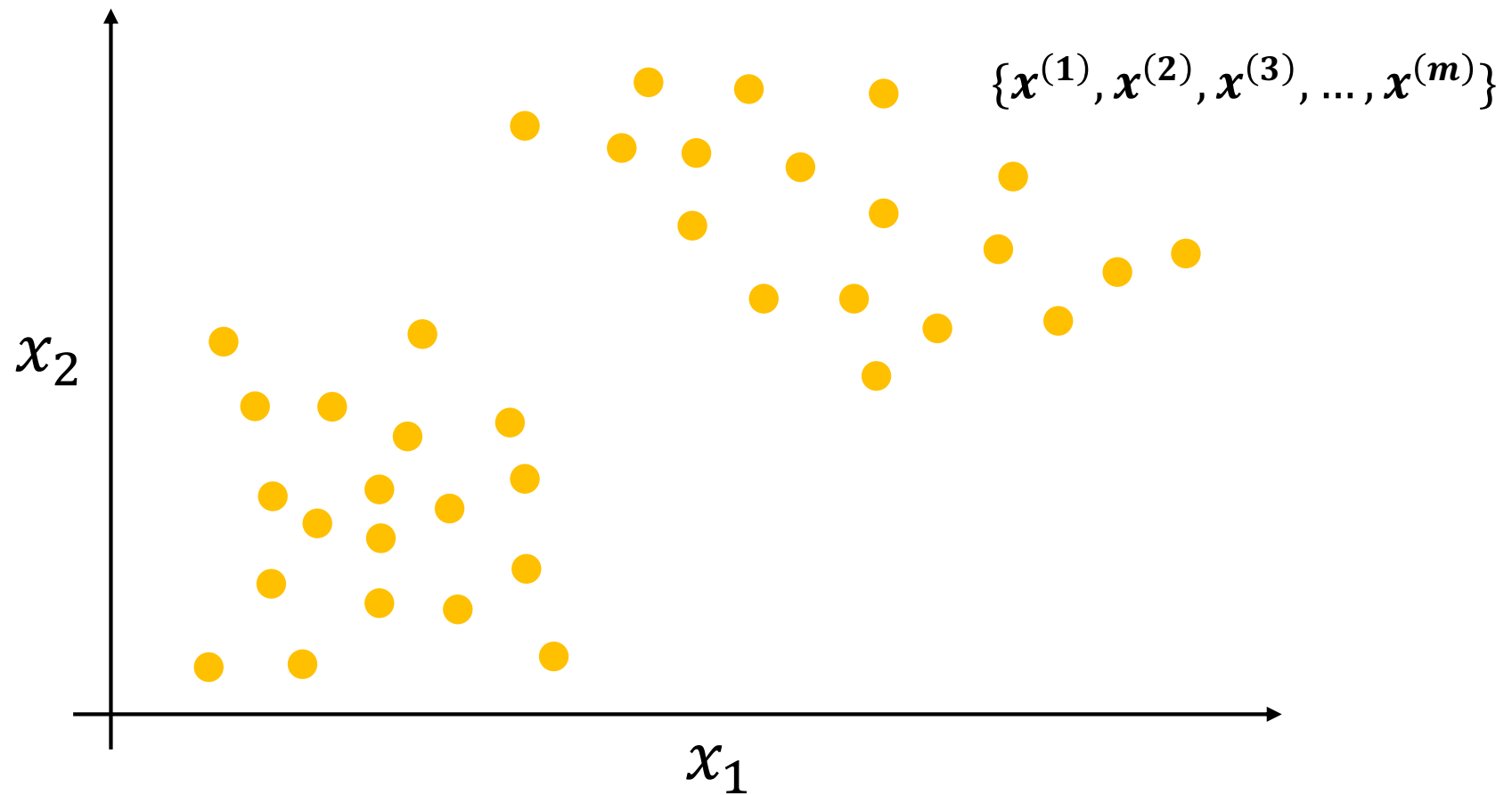
DNA analysis



Astronomical data analysis

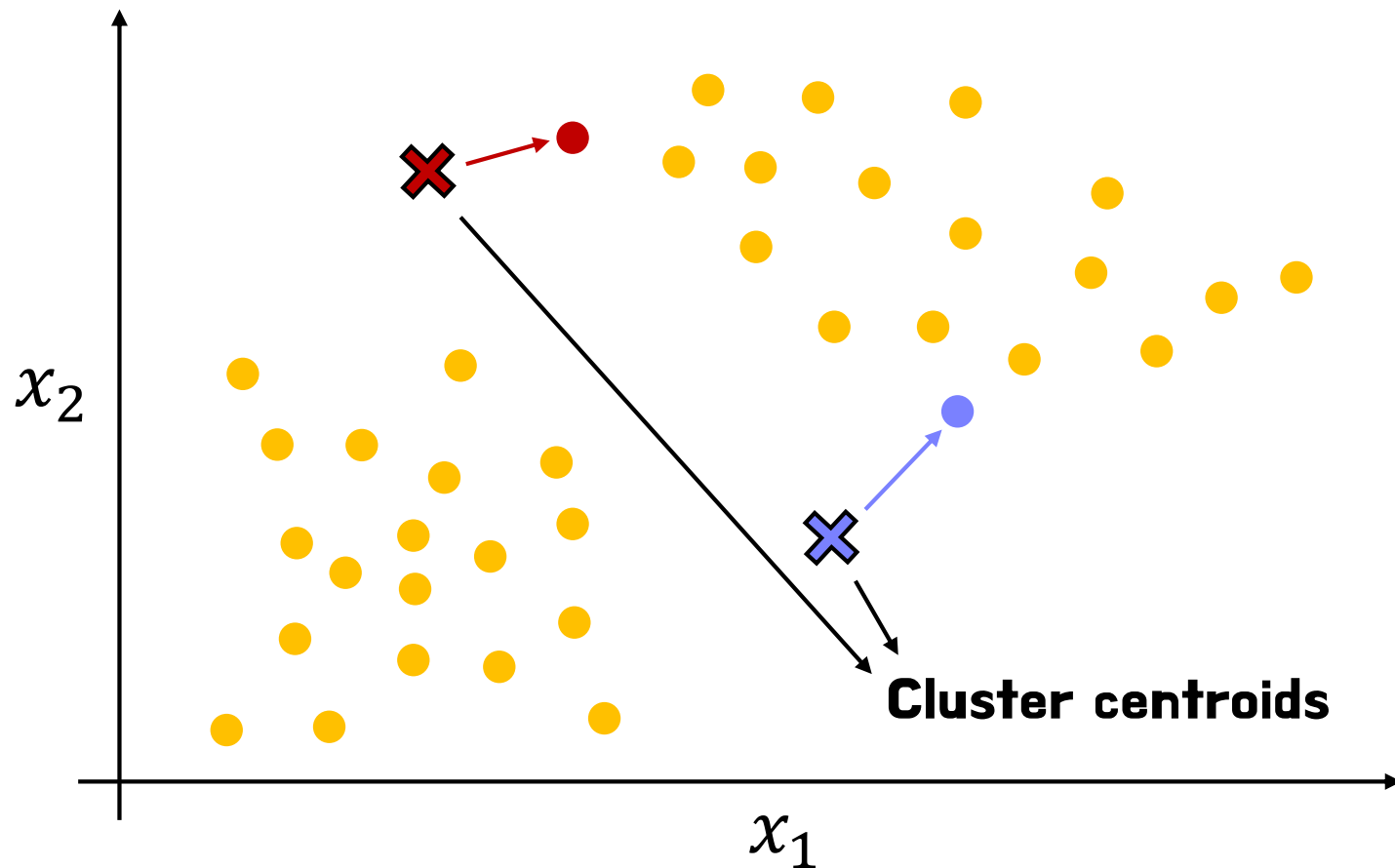


K-means intuition



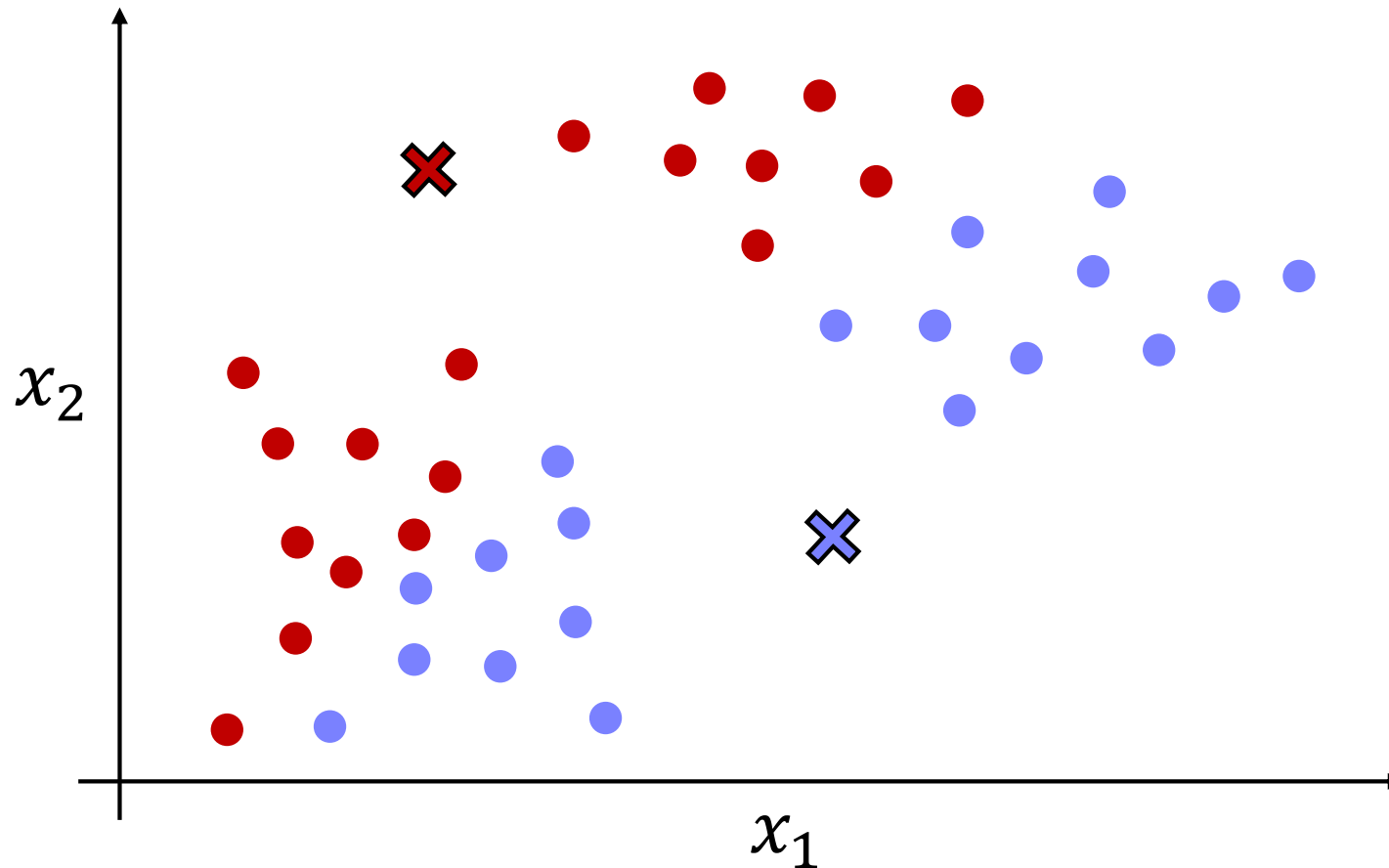
K-means intuition

Step 1: Assign each point to its closest centroid



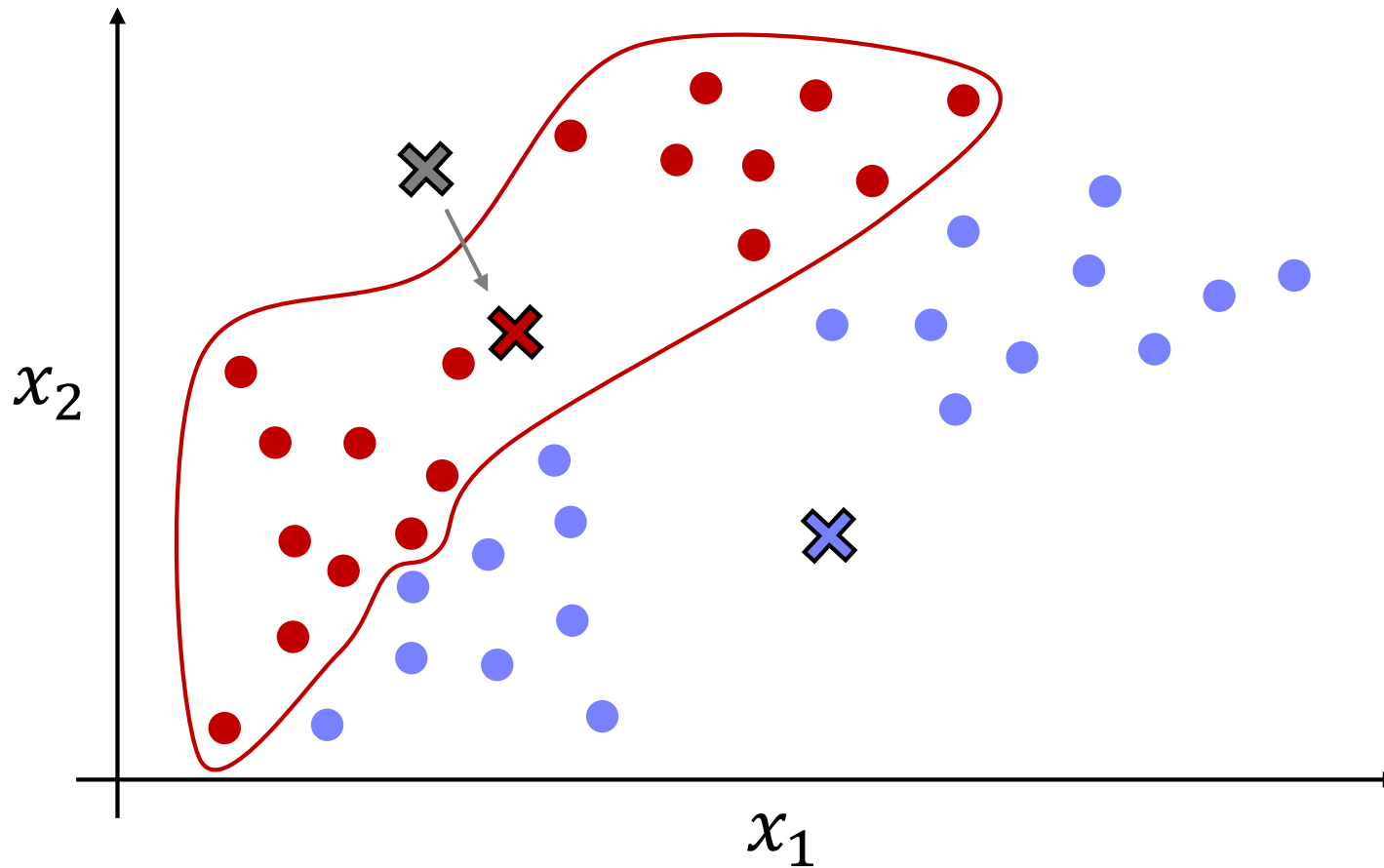
K-means intuition

Step 1: Assign each point to its closest centroid



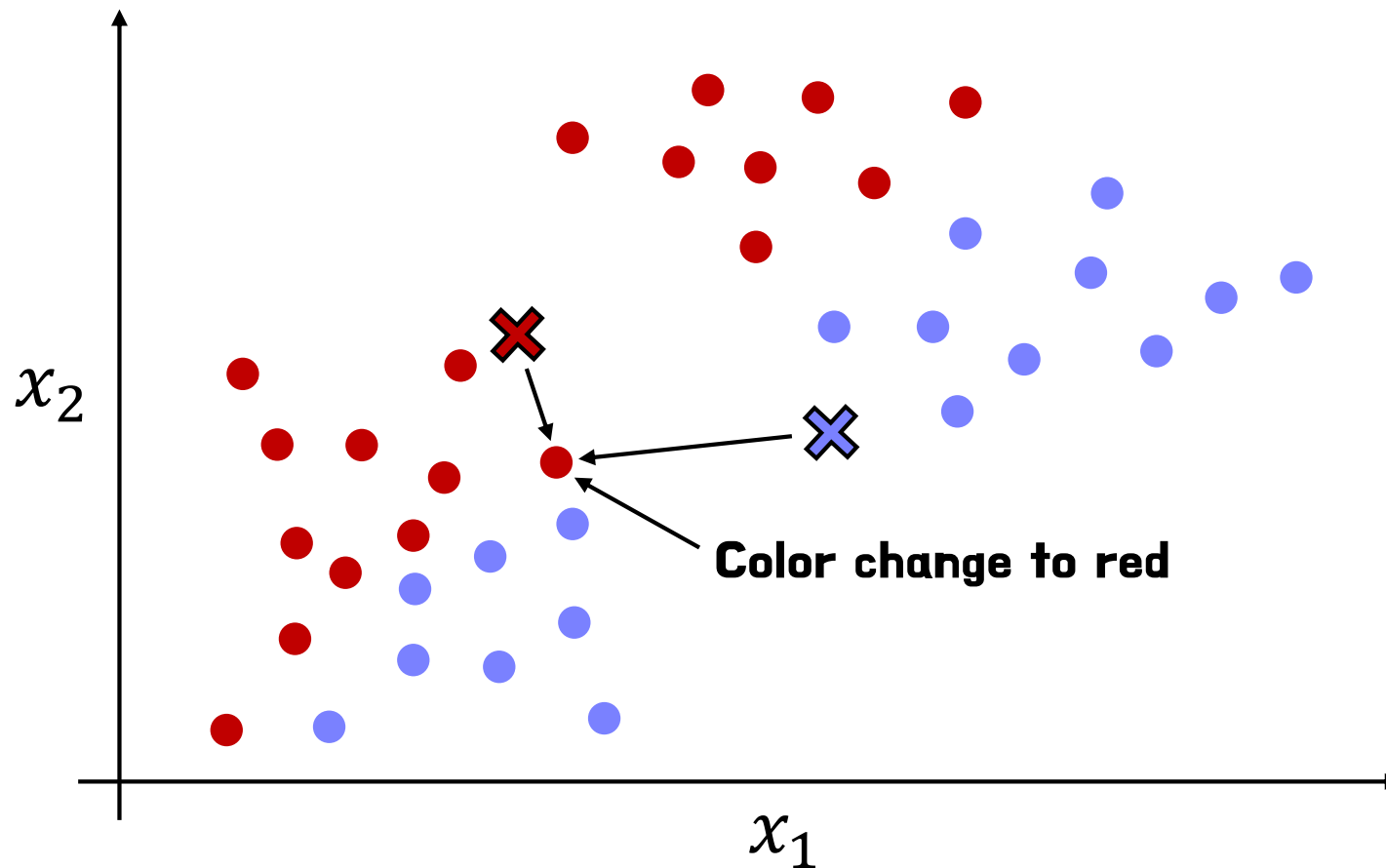
K-means intuition

Step 2: Recompute the centroids



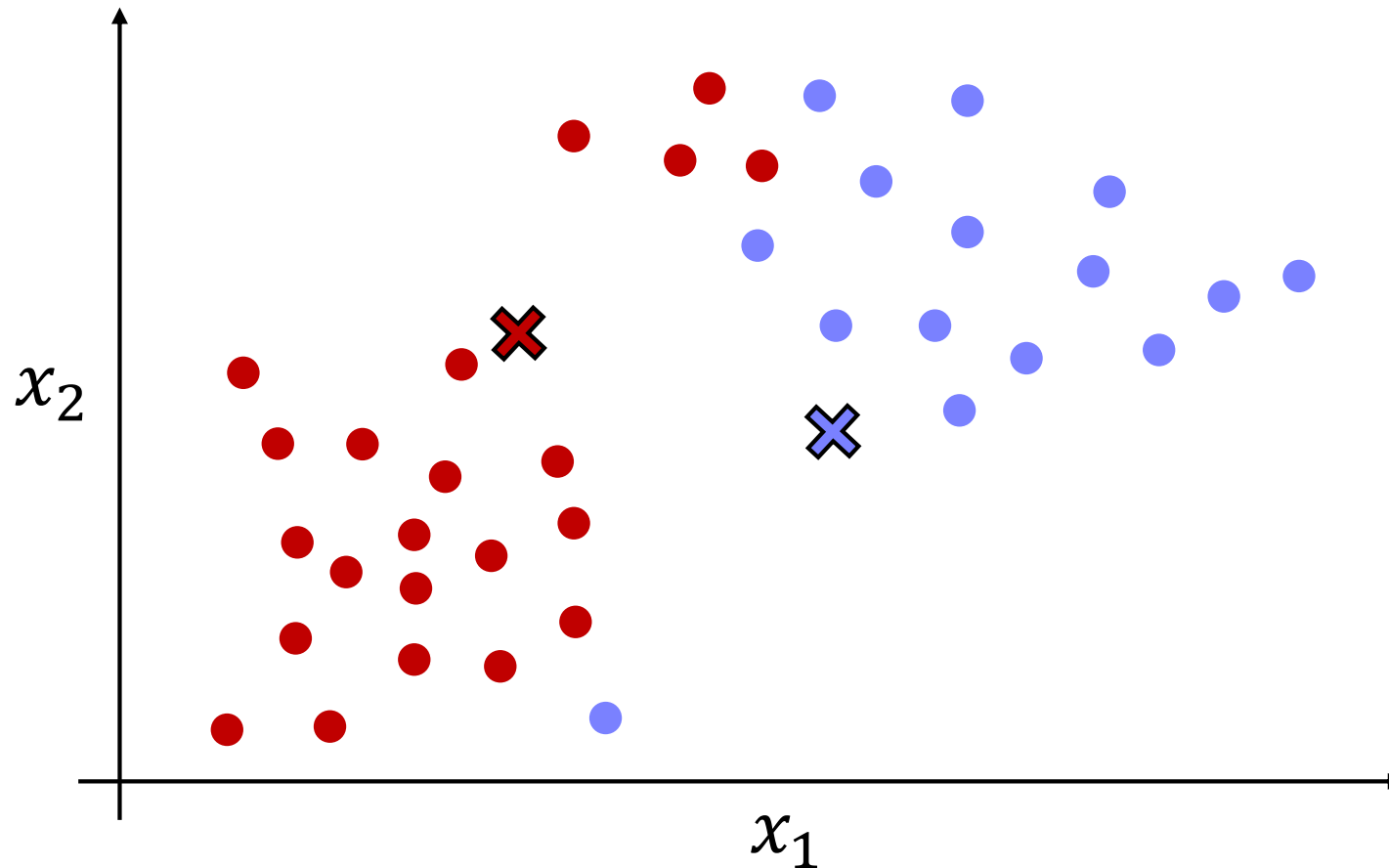
K-means intuition

Step 1: Assign each point to its closest centroid (new)



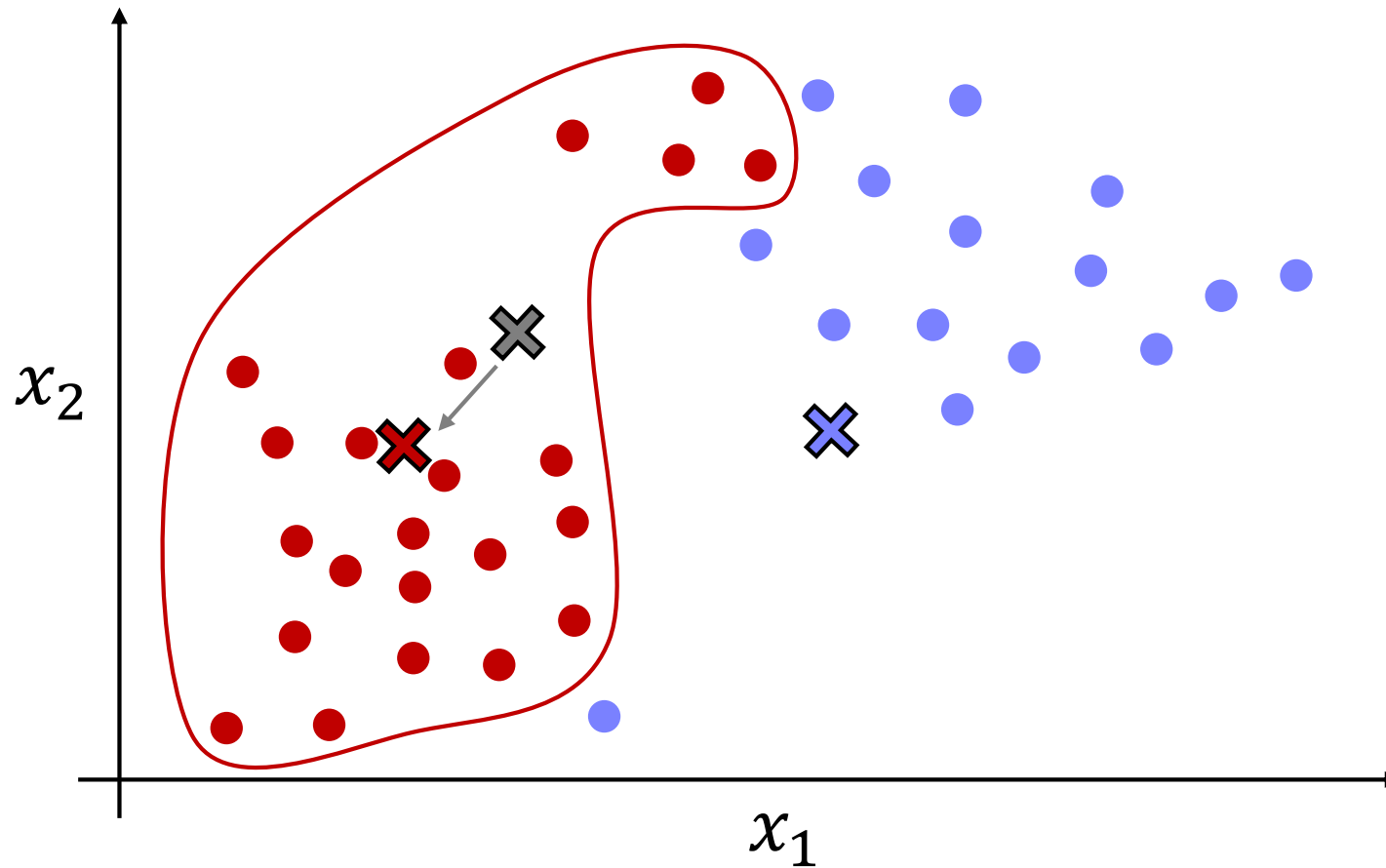
K-means intuition

Step 1: Assign each point to its closest centroid (new)



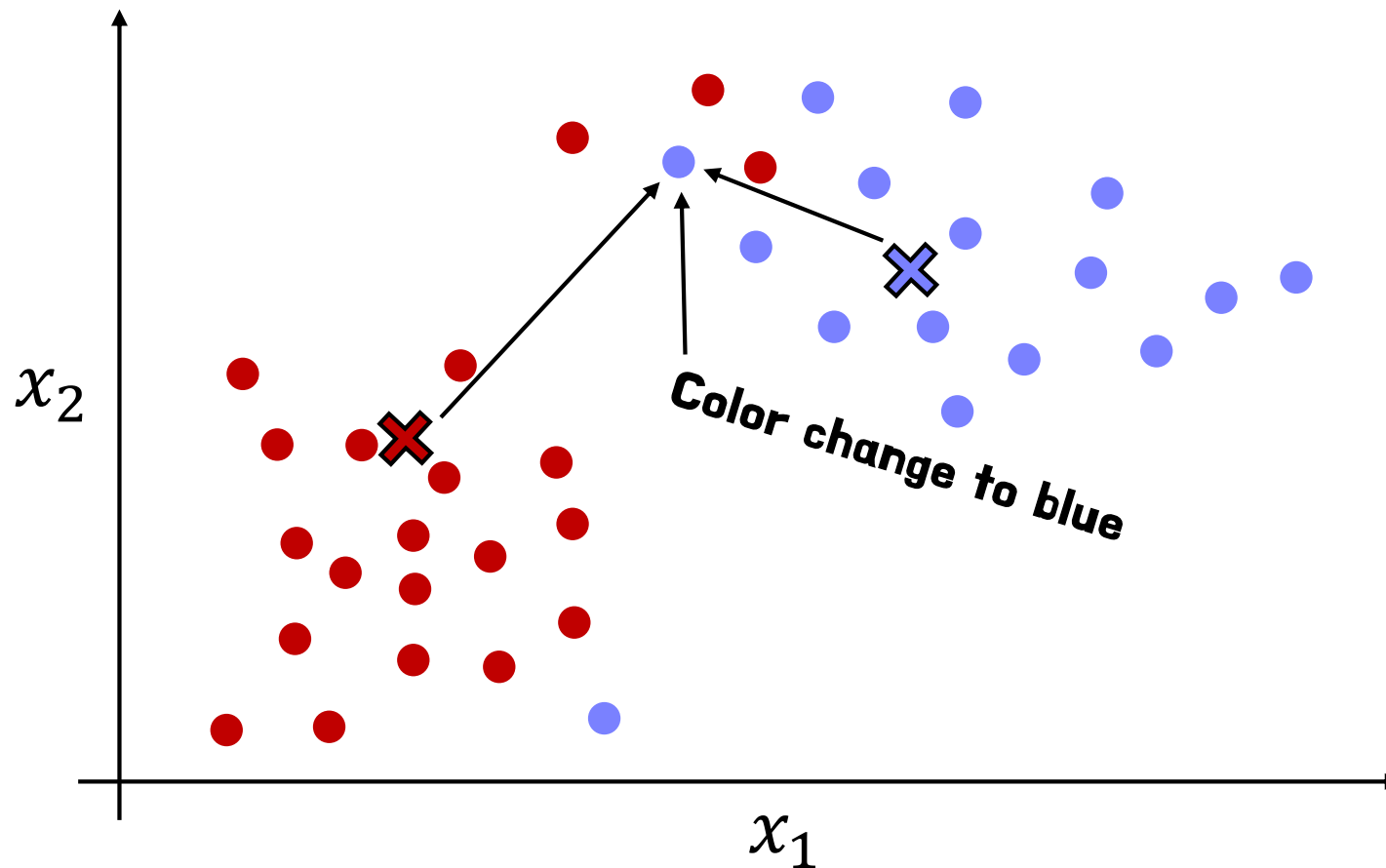
K-means intuition

Step 2: Recompute the centroids



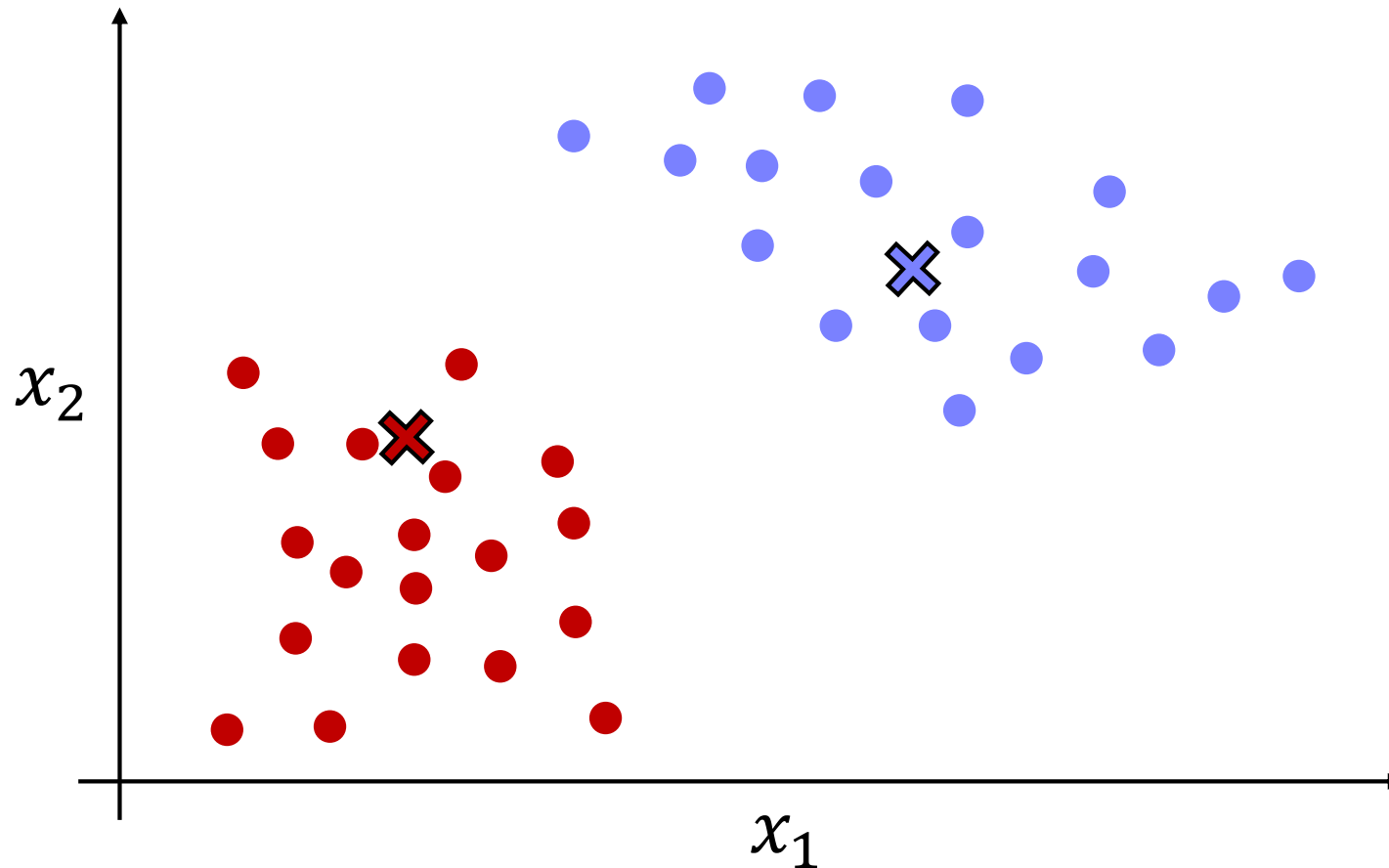
K-means intuition

Step 1: Assign each point to its closest centroid (new)



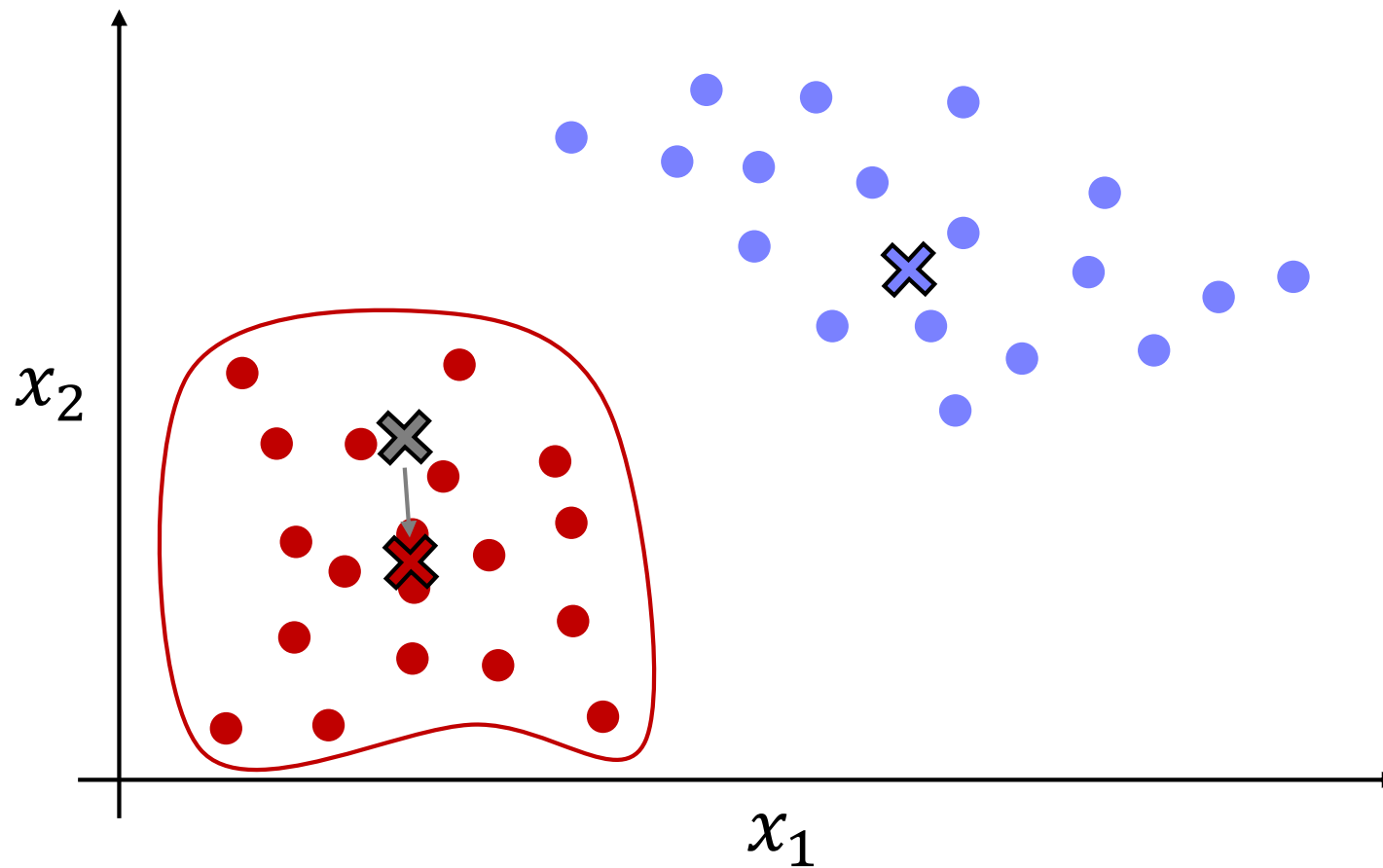
K-means intuition

Step 1: Assign each point to its closest centroid (new)



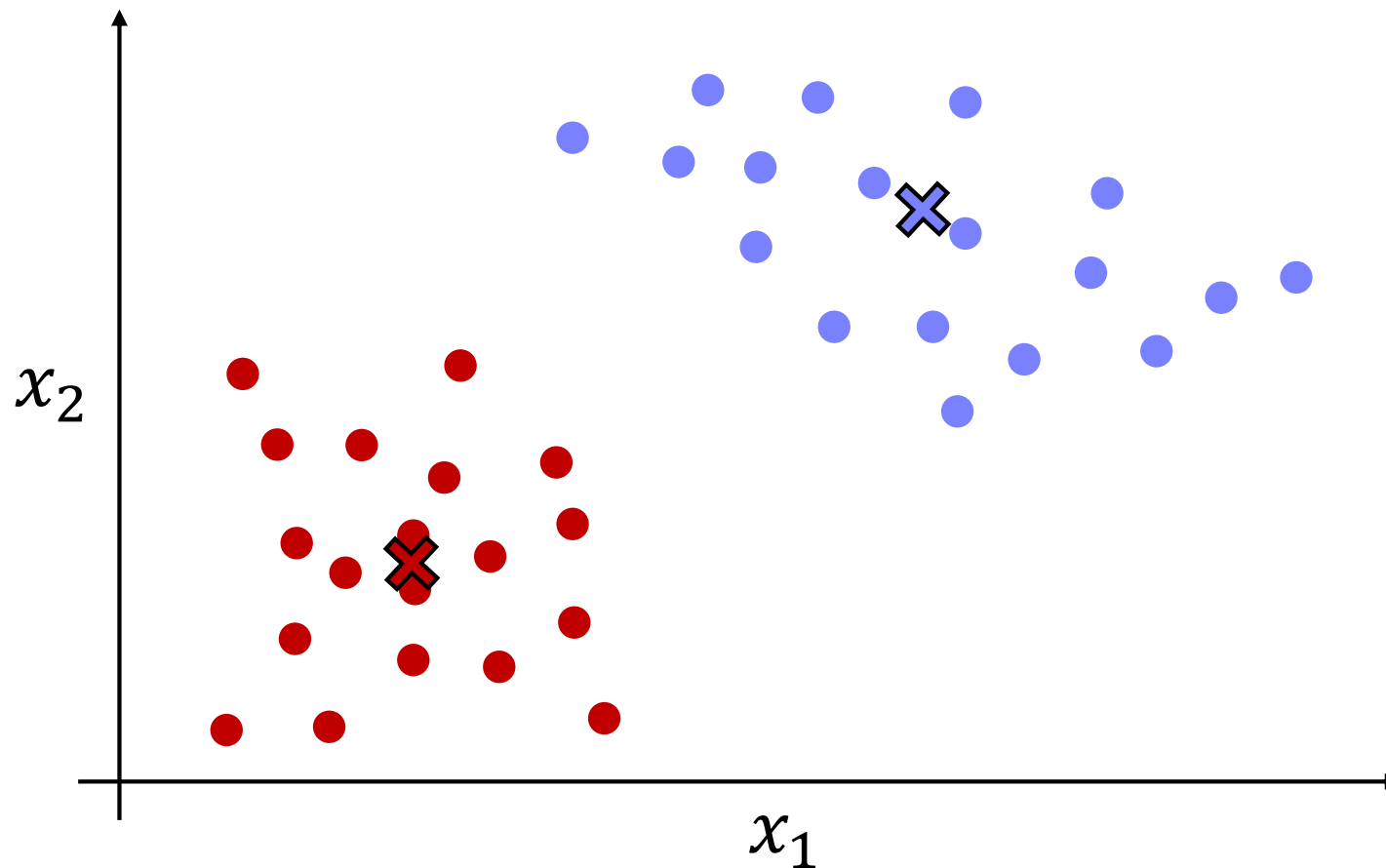
K-means intuition

Step 2: Recompute the centroids



K-means intuition

Step 1: Assign each point to its closest centroid (new)



K-means algorithm

Randomly initialize K cluster centroids $(\mu_1, \mu_2, \dots, \mu_K)$

Repeat {

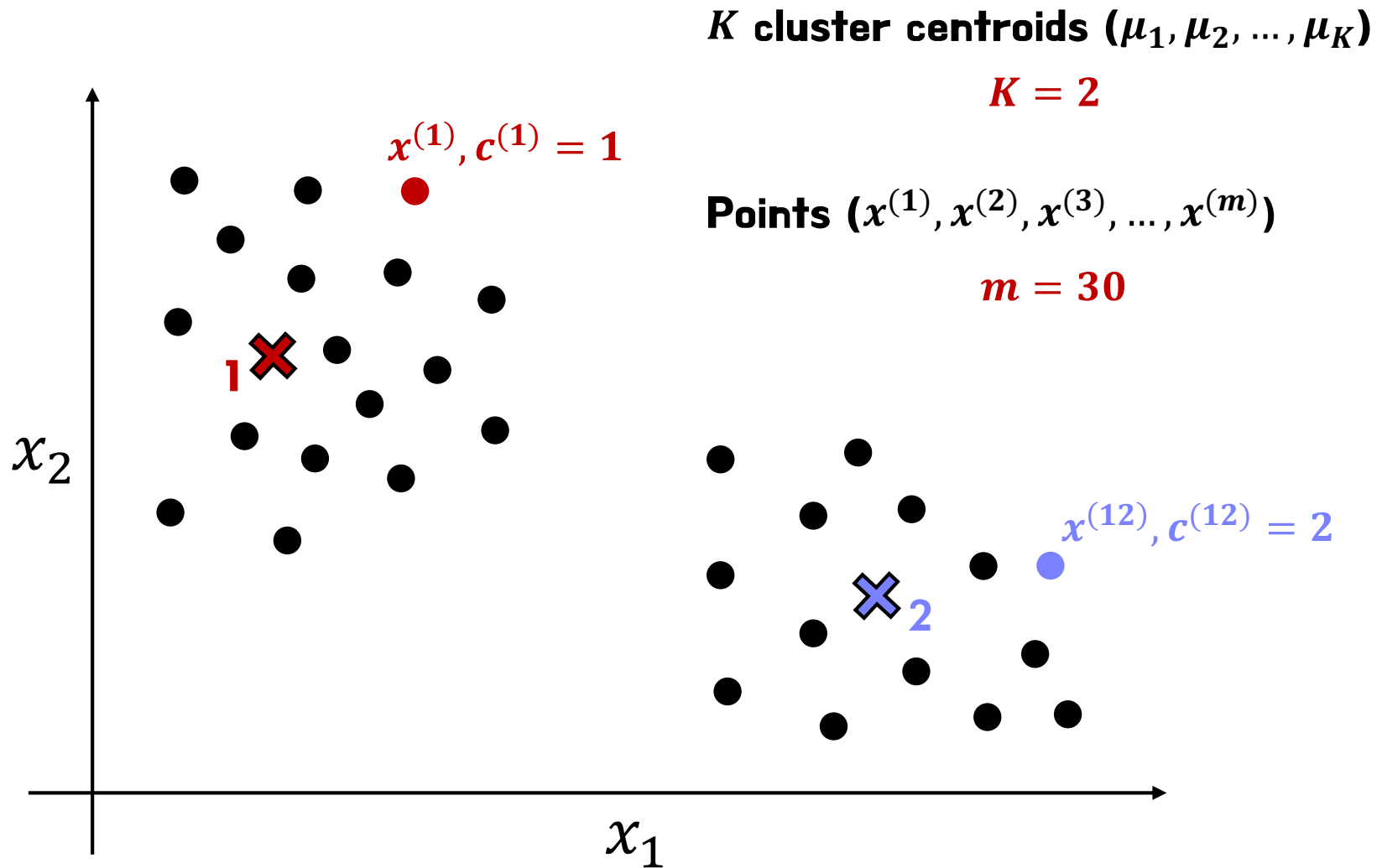
Assign points to cluster centroids

for $i = 1$ to m

$c^{(i)} :=$ index (from 1 to K) of cluster centroid closest to $x^{(i)}$

$$\begin{array}{l} \downarrow \\ \min_k \|x^{(i)} - \mu_k\|^2 \end{array}$$

K-means algorithm



K-means algorithm

Randomly initialize K cluster centroids $(\mu_1, \mu_2, \dots, \mu_K)$

Repeat {

Assign points to cluster centroids

for $i = 1$ to m

$c^{(i)} :=$ index (from 1 to K) of cluster centroid closest to $x^{(i)}$

$$\begin{array}{l} \downarrow \\ \min_k \|x^{(i)} - \mu_k\|^2 \end{array}$$

Move cluster centroids

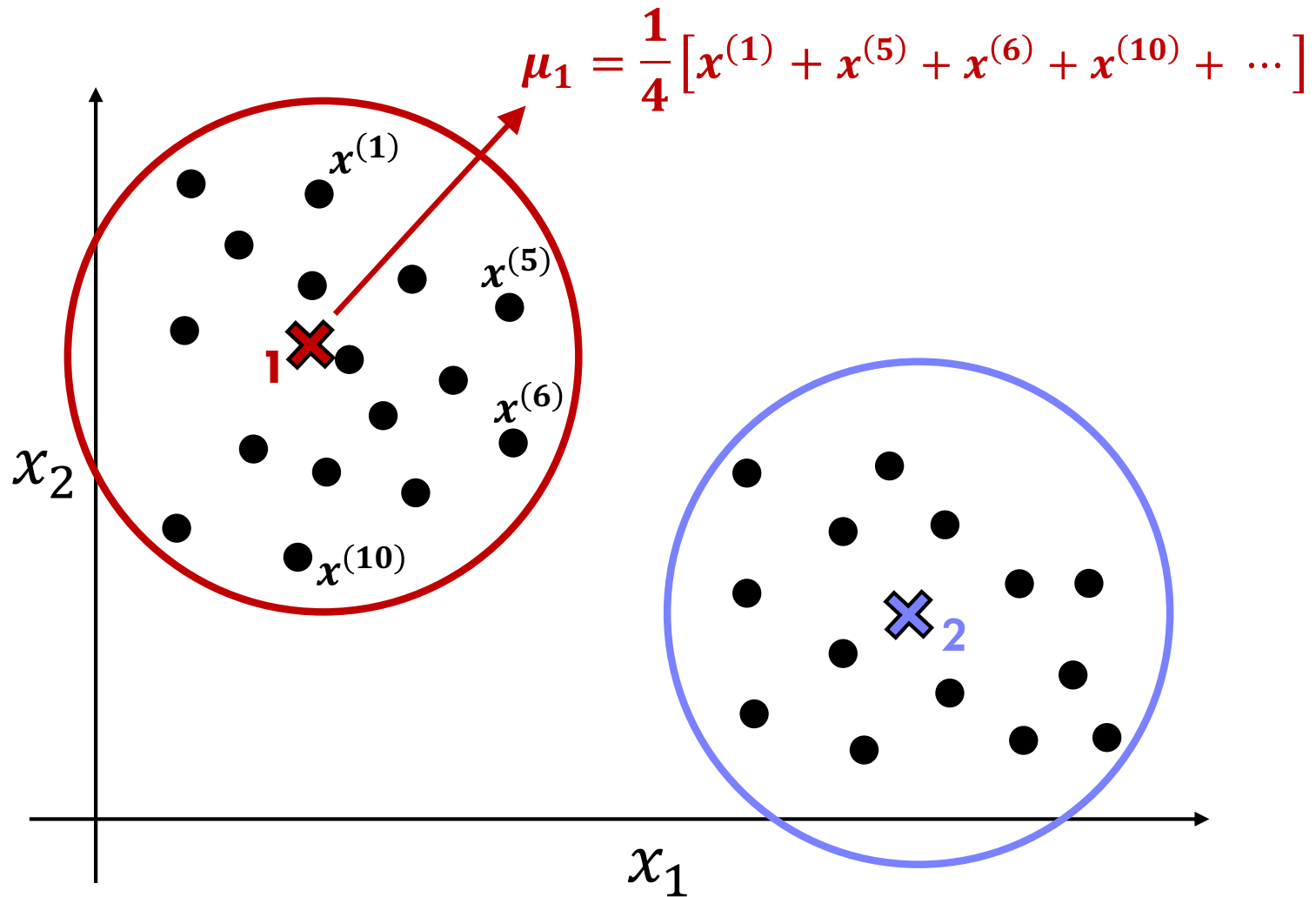
for $k = 1$ to K

$\mu_k :=$ average (mean) of points assigned to cluster k

$$\begin{array}{l} \downarrow \\ \mu_1 = \frac{1}{4} [x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)} + \dots] \end{array}$$



K-means algorithm

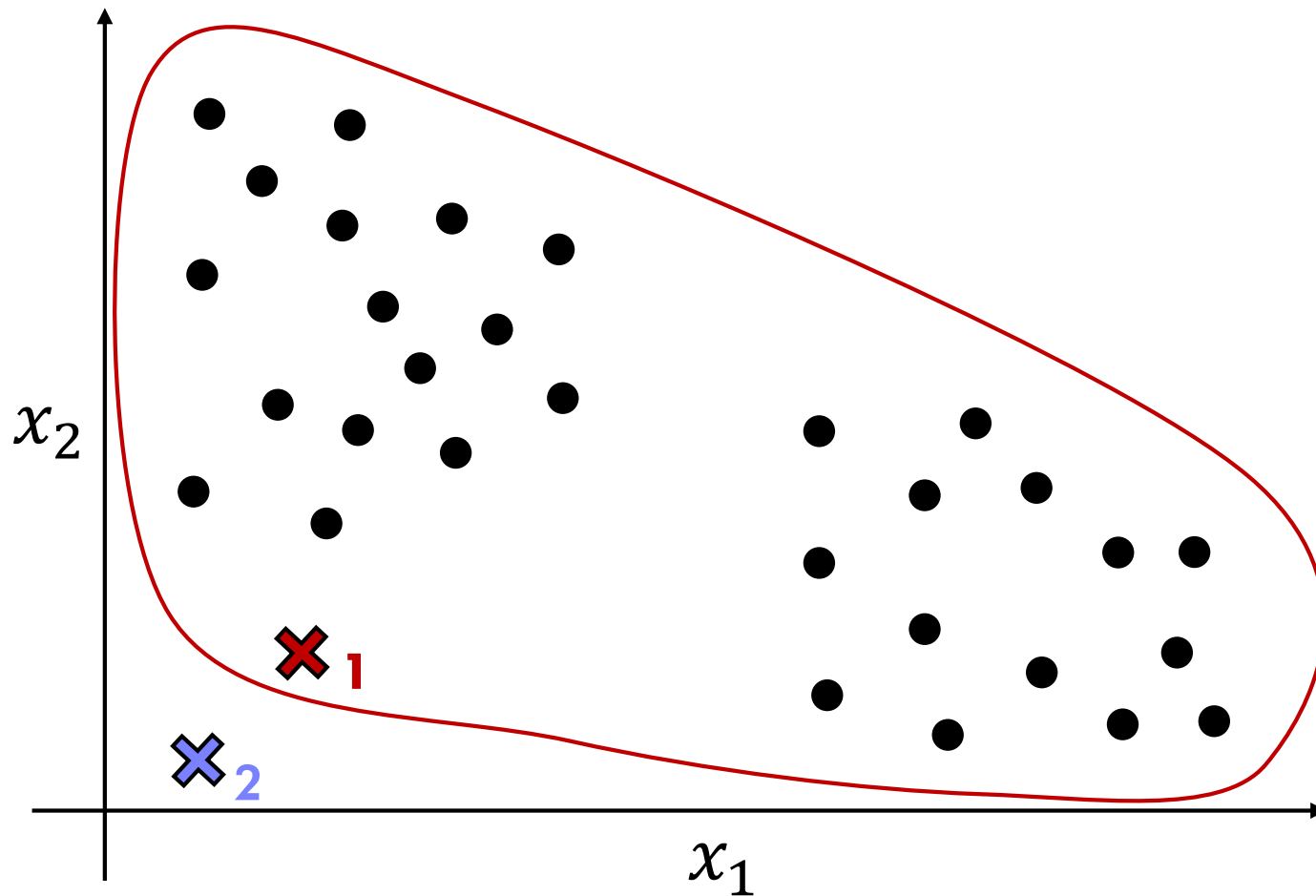


K-means algorithm

In case, one of the centroids doesn't have any point

→ simply remove one centroid ($K - 1$)

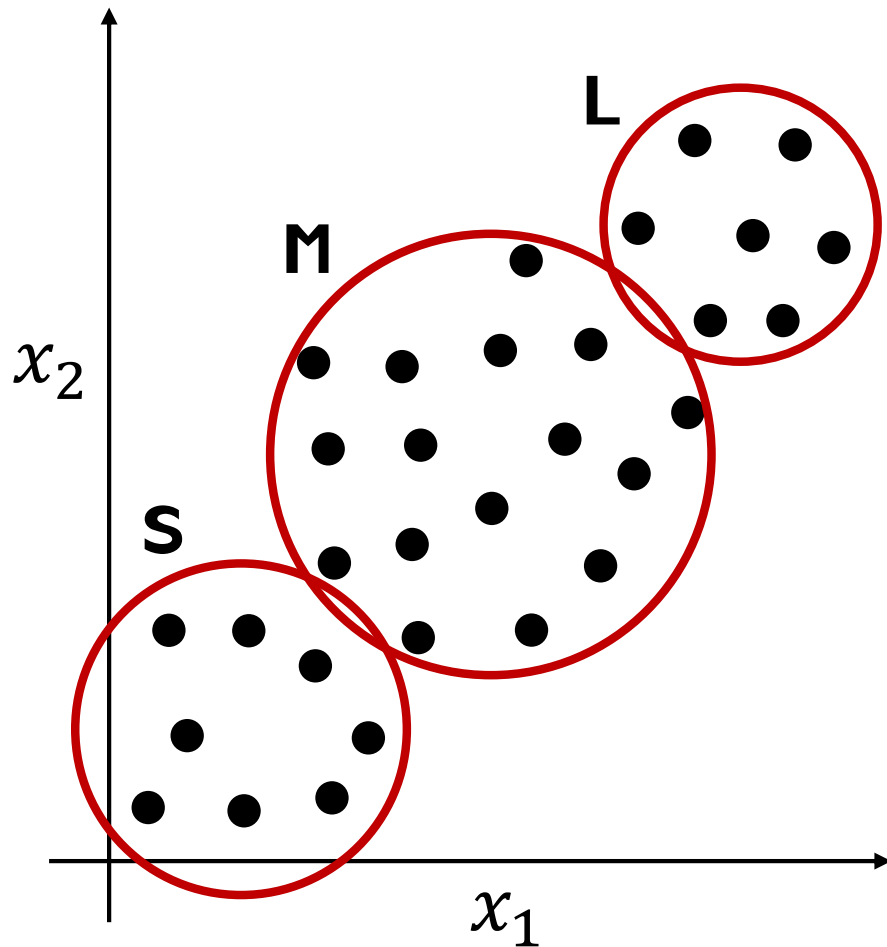
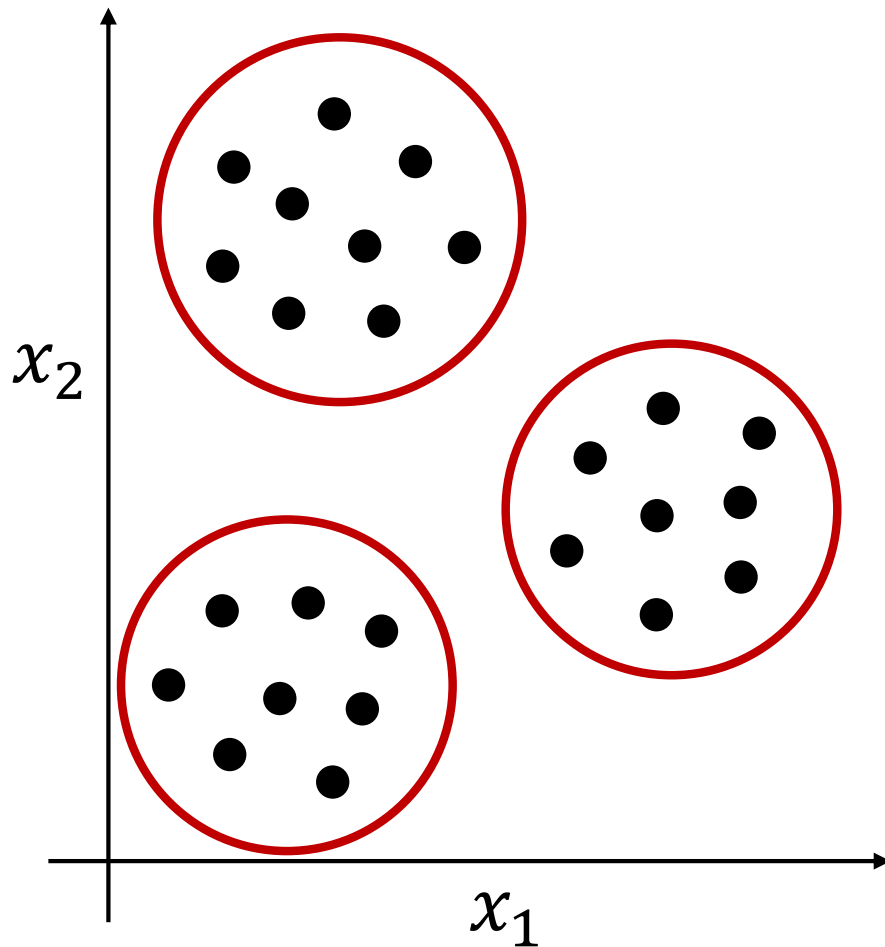
→ or regenerate them randomly



K-means algorithm

Working for points that are not well separated

T-shirt sizing



K-means optimization objective

$c^{(i)}$ = index of cluster (1,2, ..., K) to which example $x^{(i)}$ is currently assigned

μ_k = cluster centroid k

$\mu_{c^{(i)}}$ = cluster centroid of cluster to which example $x^{(i)}$ has been assigned

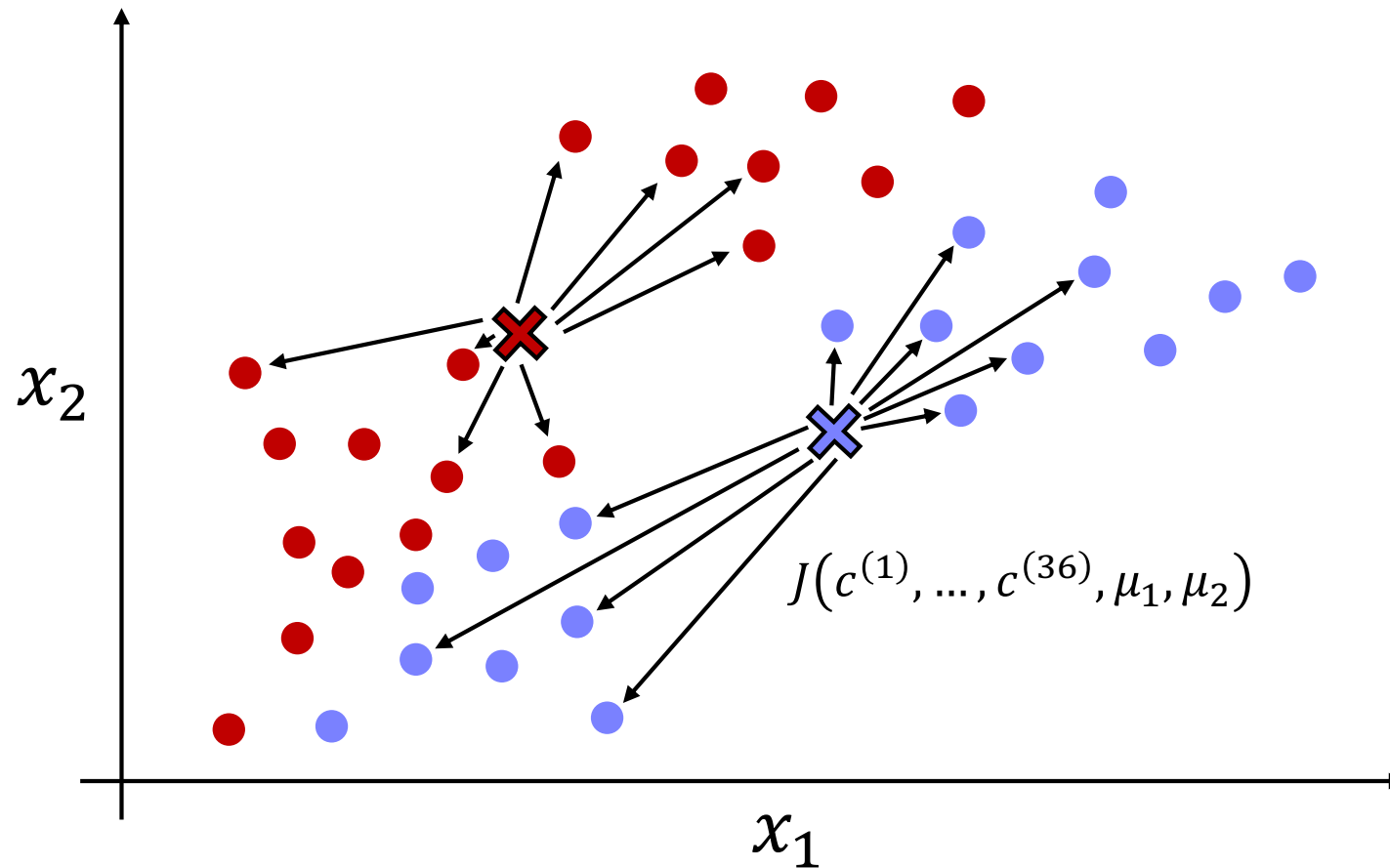
Cost function (or distortion function)

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$$\min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$$



Cost function



Cost function for K-means

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

Repeat {

Assign points to cluster centroids

for $i = 1$ to m

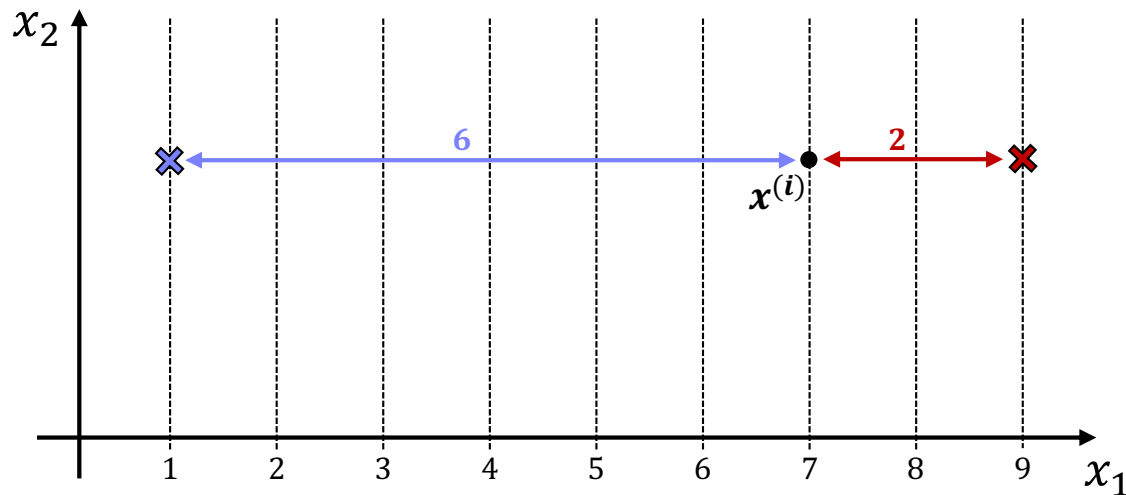
$c^{(i)} :=$ index (from 1 to K) of cluster centroid closest to $x^{(i)}$

Move cluster centroids

for $k = 1$ to K

$\mu_k :=$ average (mean) of points assigned to cluster k

}



Cost function for K-means

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

Repeat {

Assign points to cluster centroids

for $i = 1$ to m

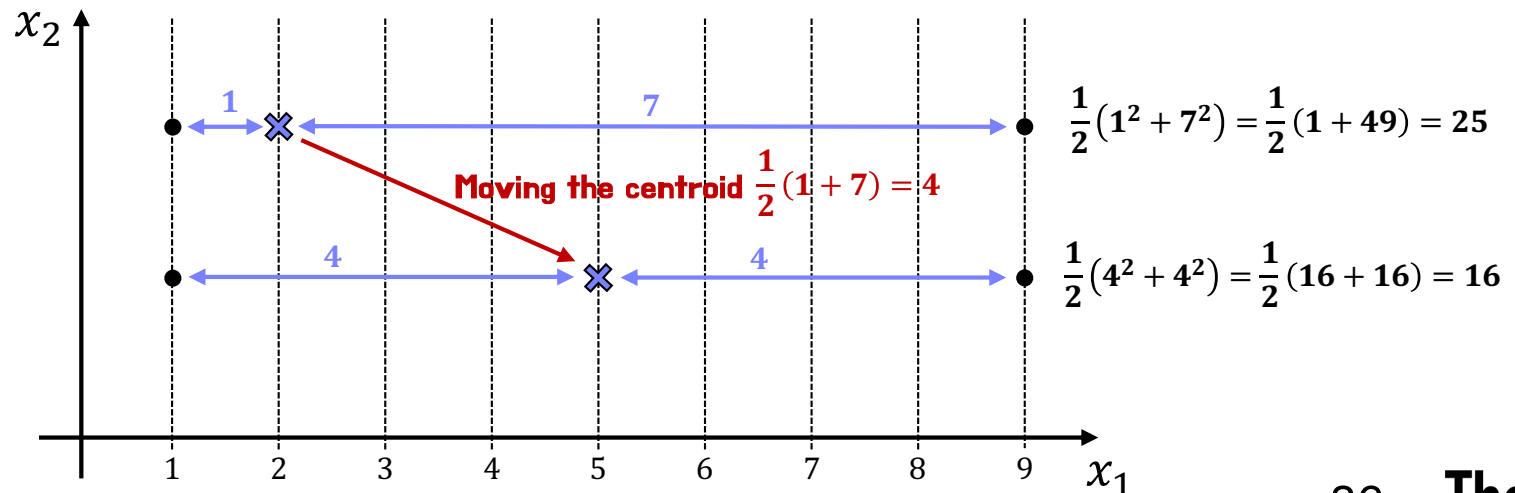
$c^{(i)} :=$ index (from 1 to K) of cluster centroid closest to $x^{(i)}$

Move cluster centroids

for $k = 1$ to K

$\mu_k :=$ average (mean) of points assigned to cluster k

}



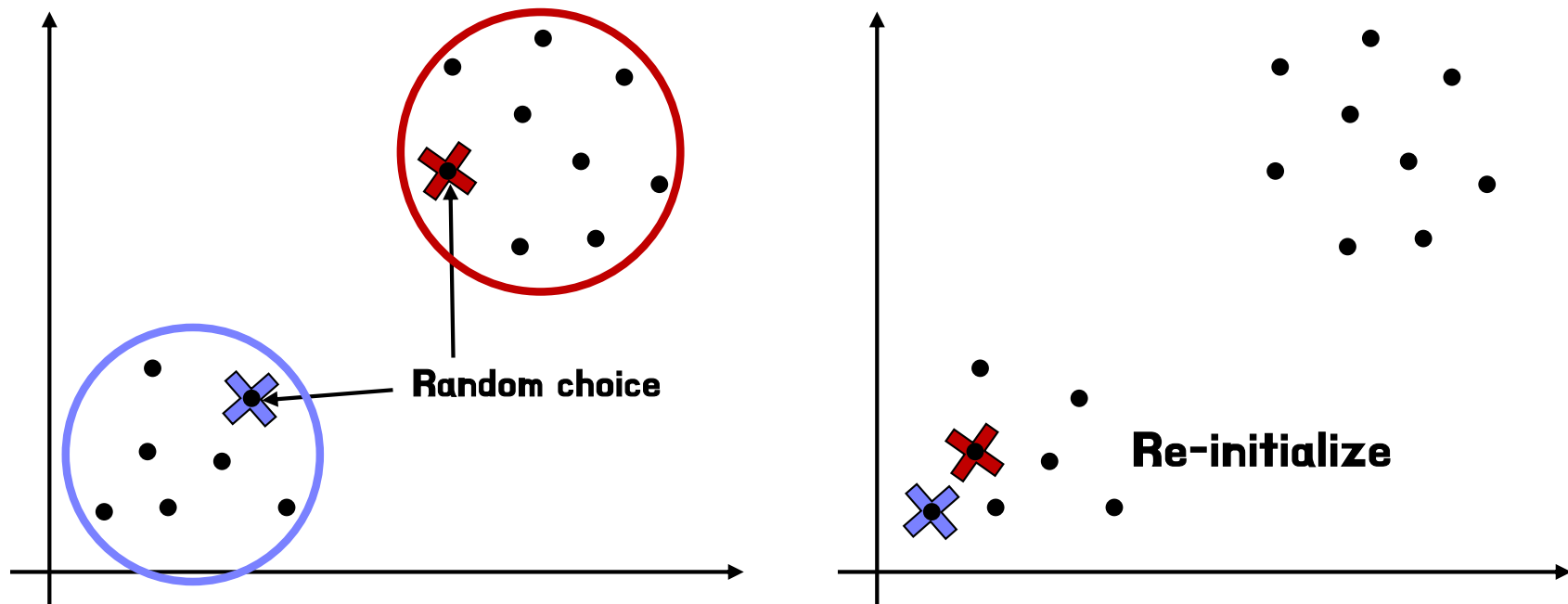
Random initialization

Choose $K < m$

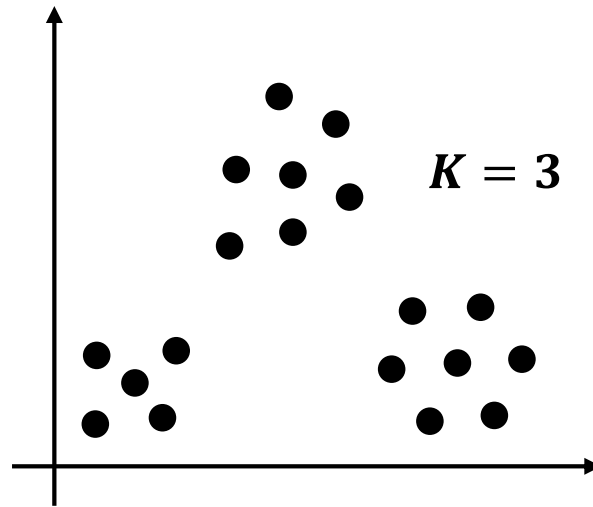
Randomly pick K training examples

Set $\mu_1, \mu_2, \dots, \mu_k$ equal to these K examples

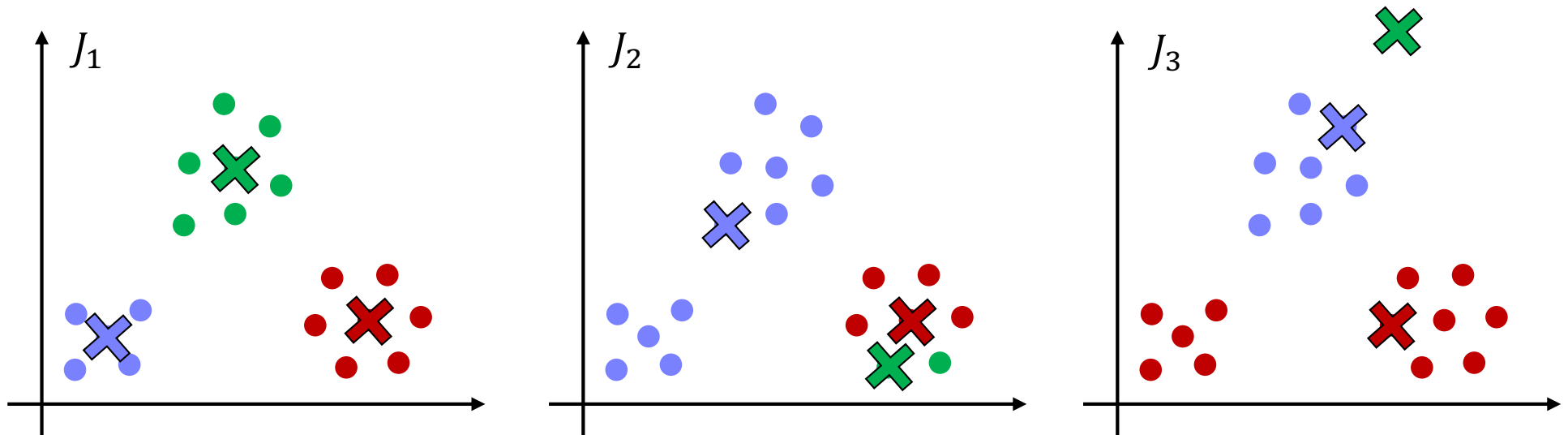
Common way to initialize the K-means optimization



Random initialization



Compare the amount of cost function



Random initialization

50 – 1,000 (commonly used)

For $i = 1$ to **100** {

Randomly initialize K-means.

Run K-means. Get $c^{(i)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k$

Computer cost function (distortion)

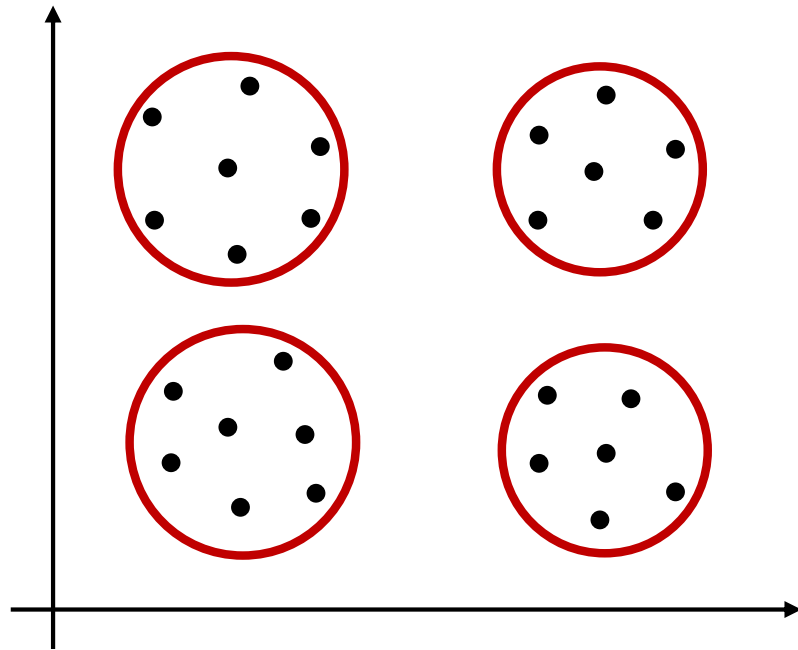
$J(c^{(i)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$

}

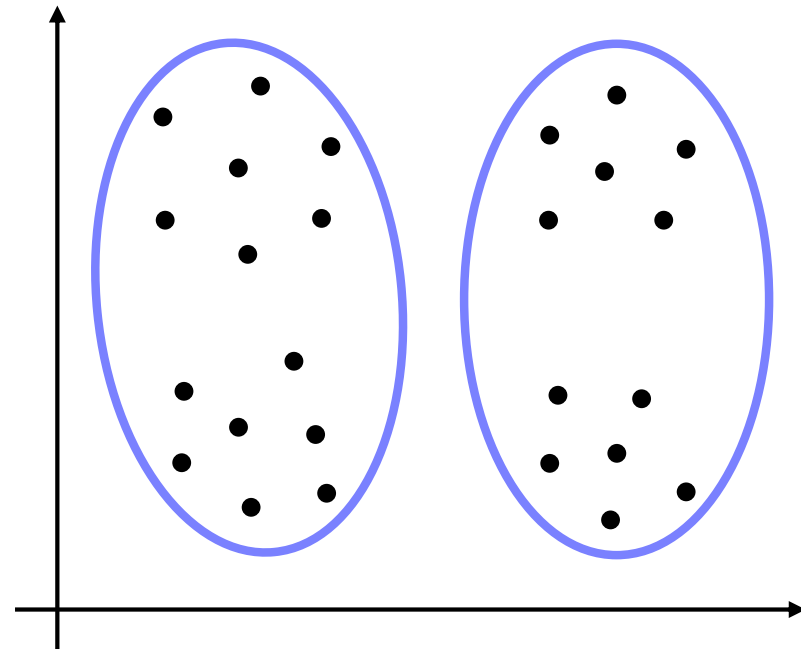
Pick set of clusters that gave lowest cost J

What is the right value of K ?

$K = 4$



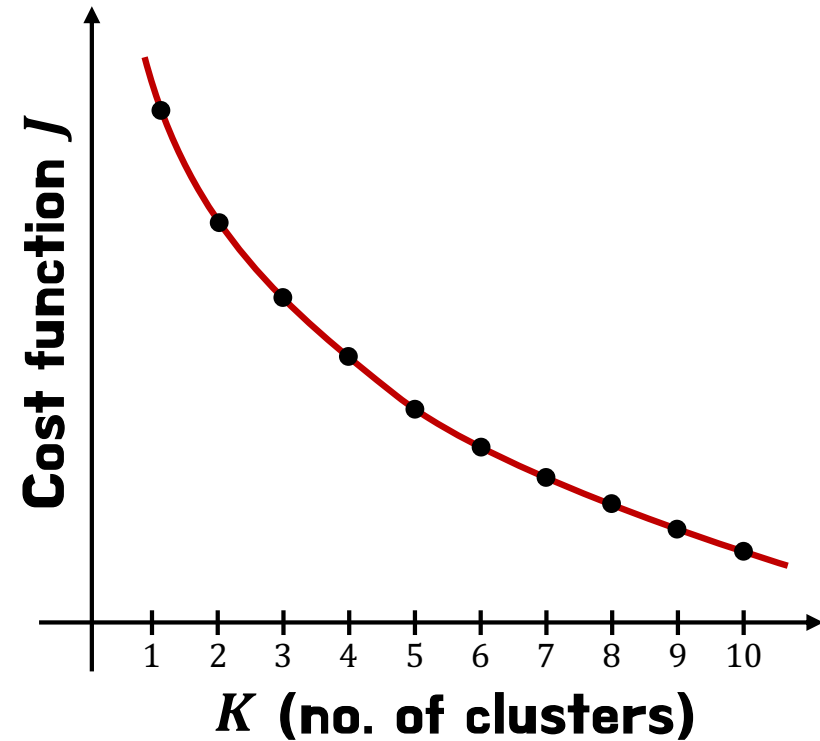
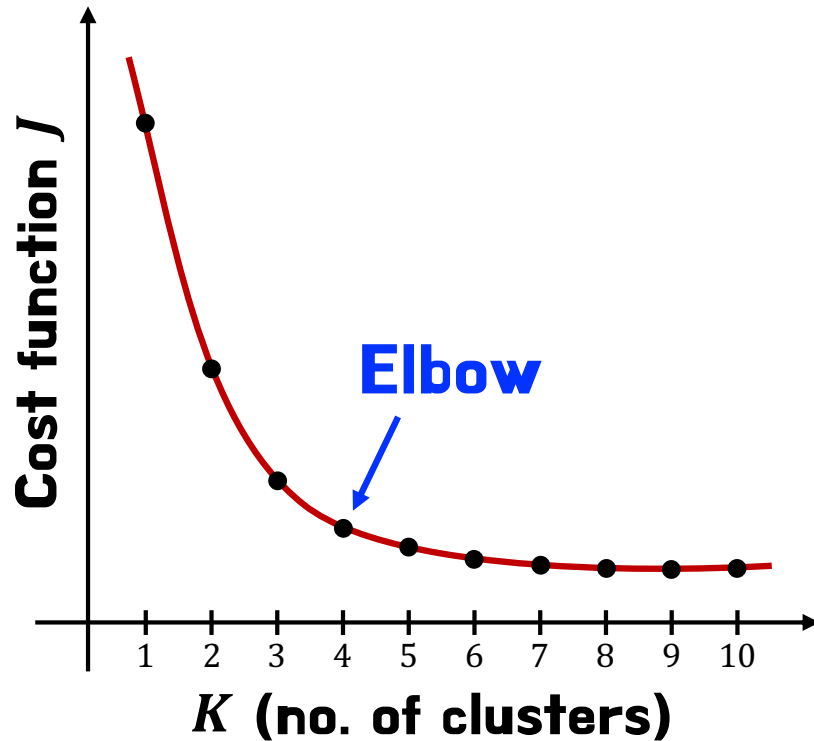
$K = 2$



Choosing the value of K

Elbow method - usually not recommended

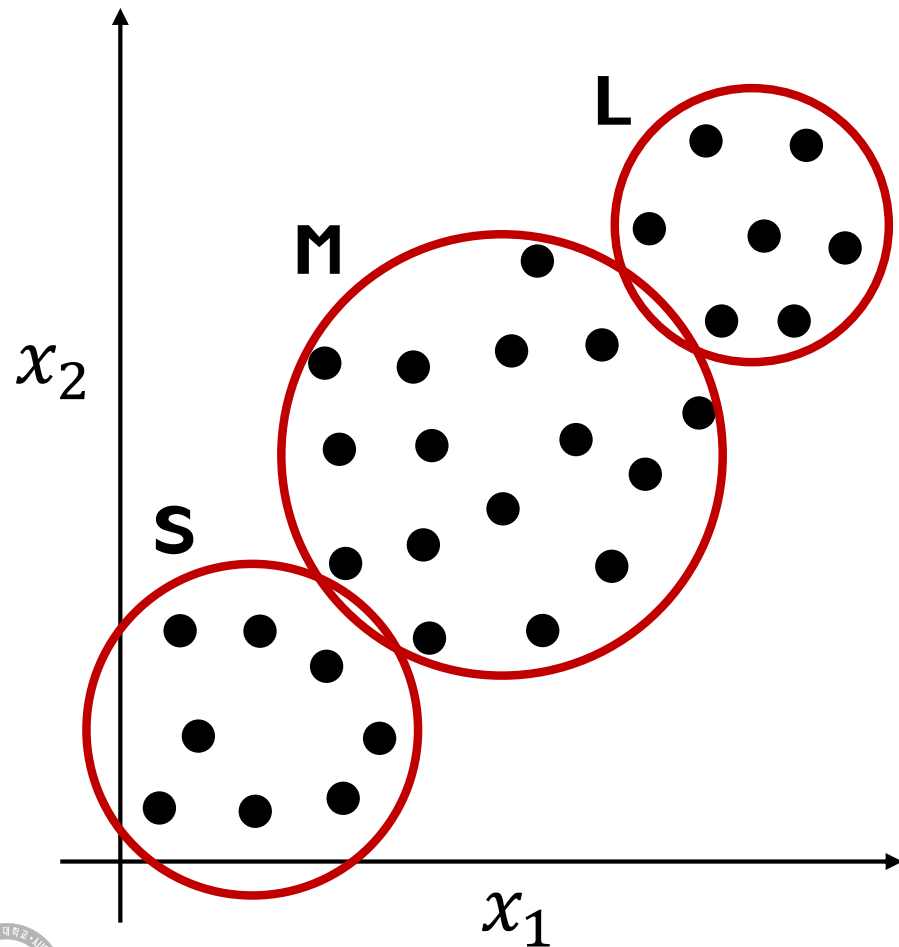
**The right " K " is often ambiguous
Don't choose K just to minimize cost J**



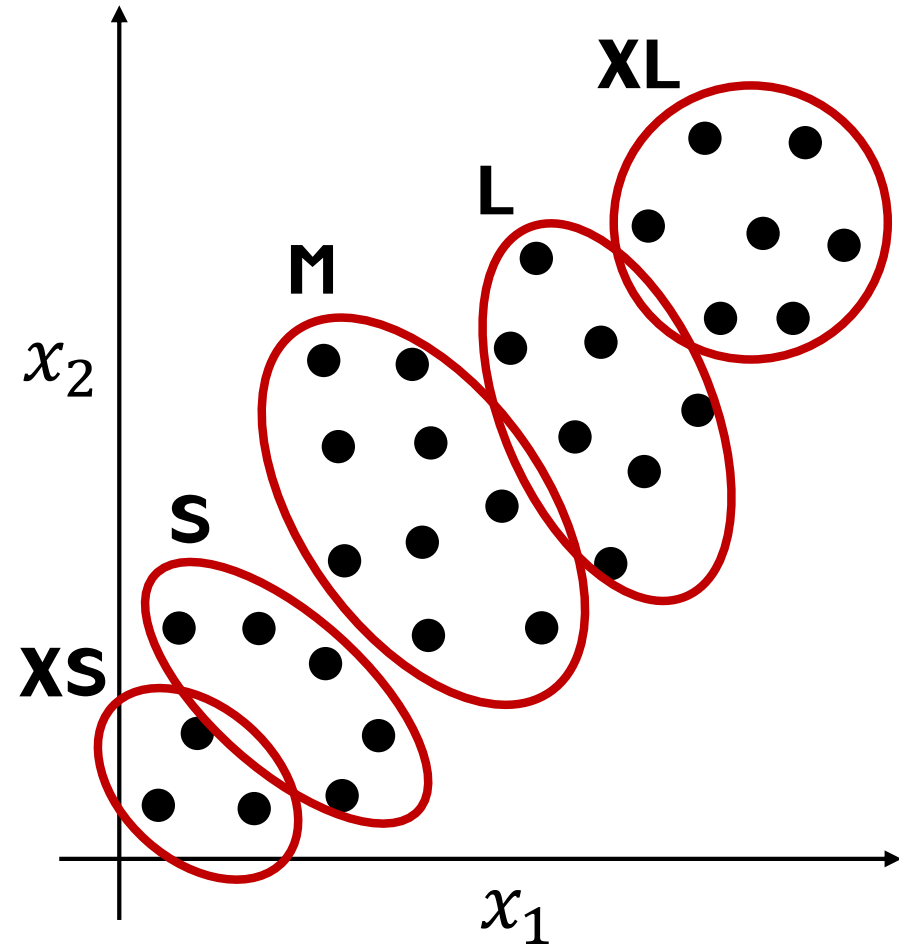
Choosing the value of K

Can be determined by later purpose

T-shirt sizing



T-shirt sizing



Anomaly detection



Anomaly detection example

Aircraft engine features:

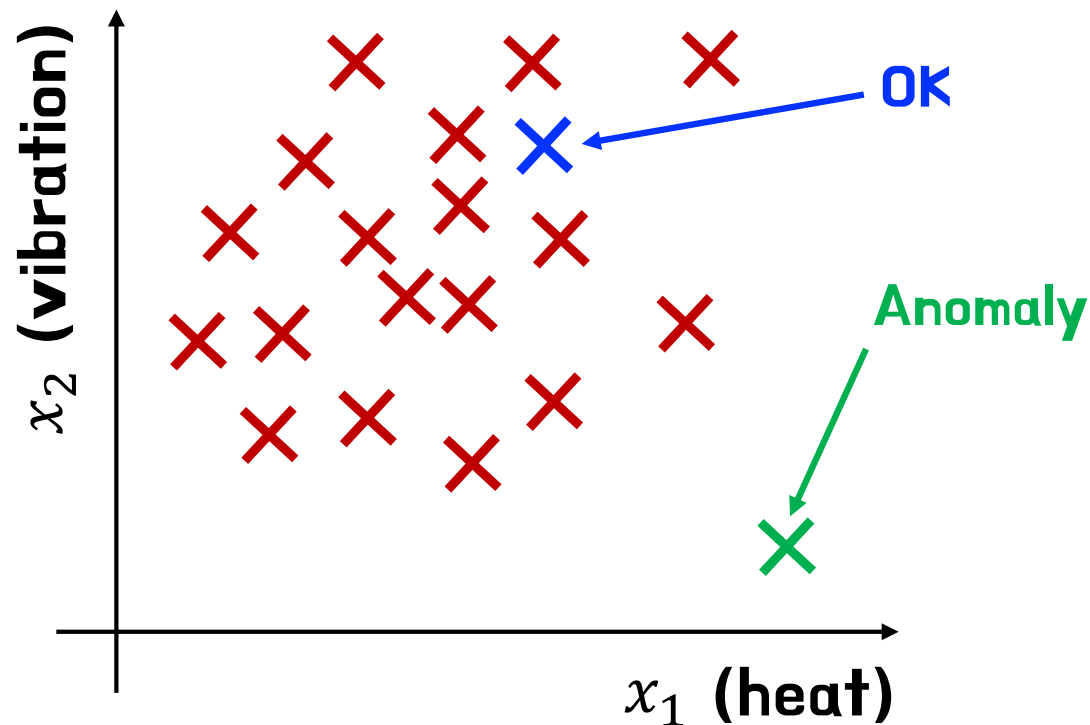
x_1 = heat generated

x_2 = vibration intensity

...

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

New engine: x_{test}



Density estimation

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

Probability of x being seen in dataset

Model $p(x)$

Is x_{test} anomalous ?

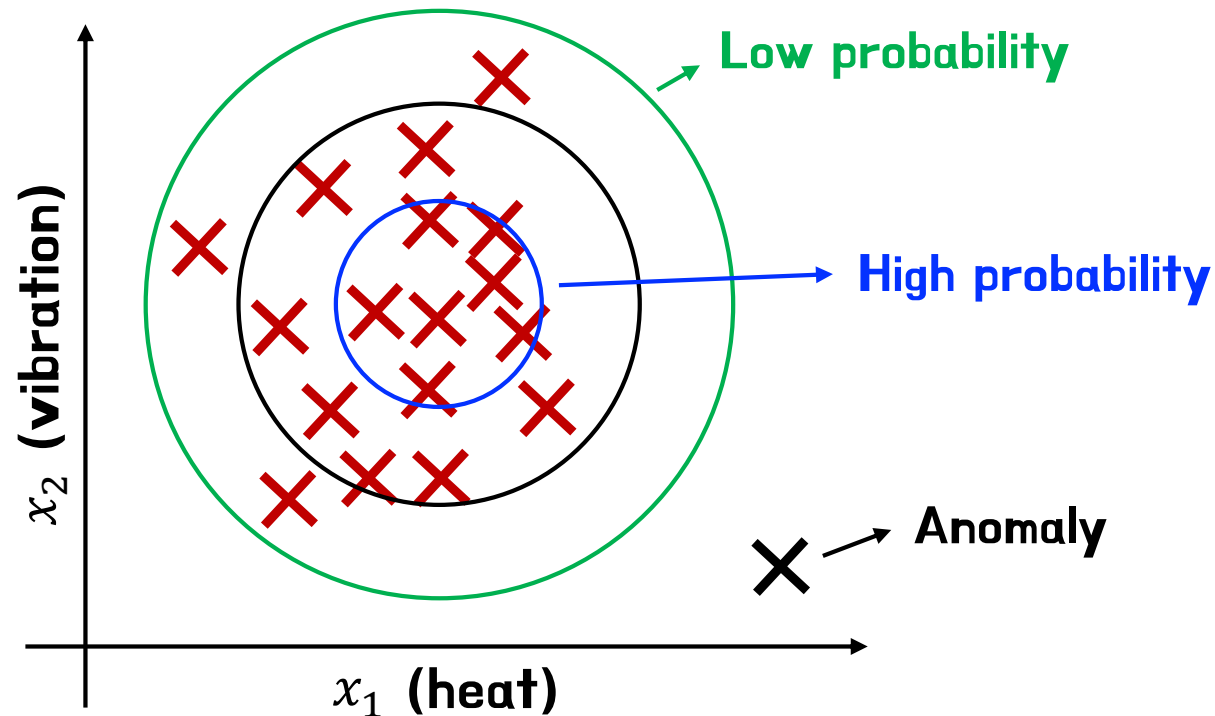
epsilon: small number

$$p(x_{test}) \geq \epsilon$$

: **OK (normal)**

$$p(x_{test}) < \epsilon$$

: **Anomaly**



Anomaly detection example

Fraud detection:

- $x^{(i)}$ = features of user i 's activities
 - Model $p(x)$ from data.
 - Identify unusual users by checking which have $p(x) < \varepsilon$
- How often log in ?
 - How many webpages visited ?
 - Transactions ?
 - Posts ?
 - Typing speed ?

Manufacturing:

- $x^{(i)}$ = features of product i
 - Airplane engine
 - Circuit board
 - Smartphone

Monitoring computers in a data center:

- $x^{(i)}$ = features of product i
 - x_1 = memory use,
 - x_2 = number of disk accesses/sec,
 - x_3 = CPU load,
 - x_4 = CPU load/network traffic,

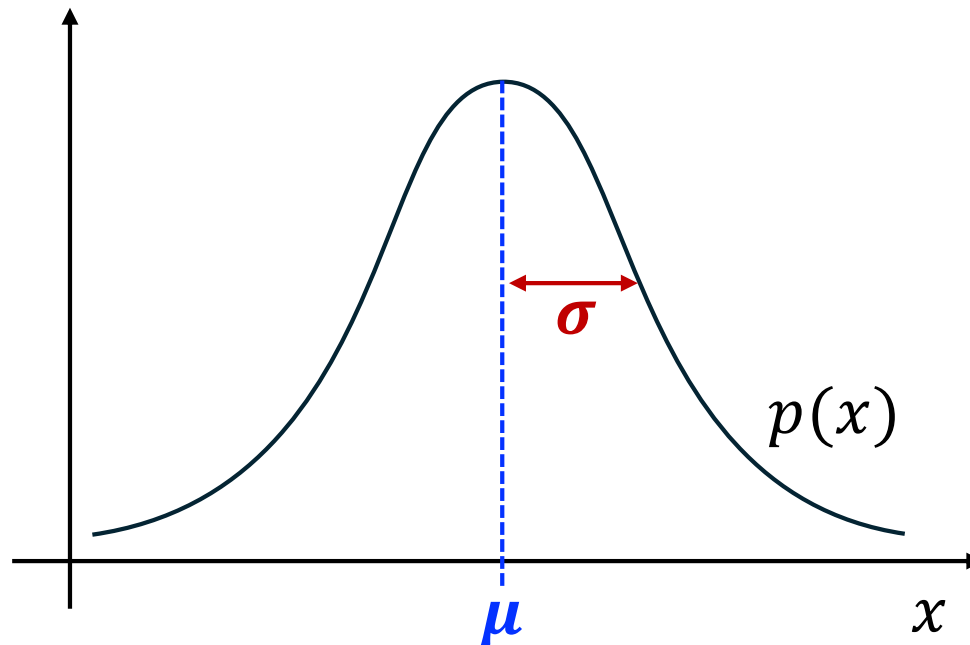


Gaussian (Normal) distribution

Say x is a number.

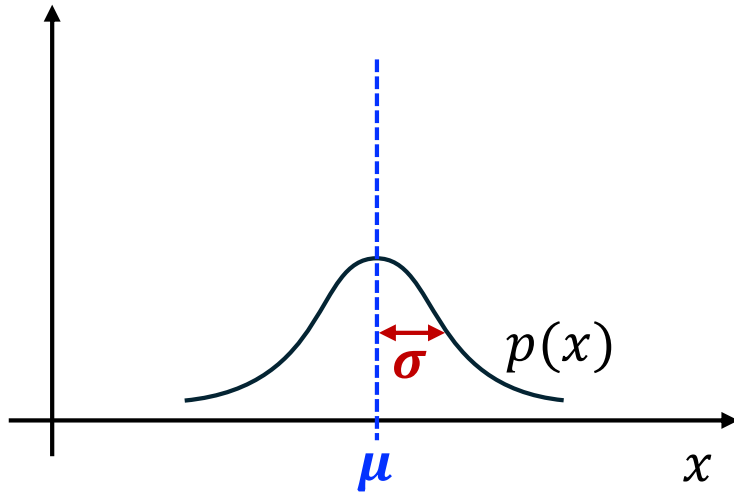
Probability of x is determined by a Gaussian with mean μ , variance σ^2 .

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad \begin{array}{l} \sigma : \text{standard deviation} \\ \sigma^2 : \text{variance} \end{array}$$

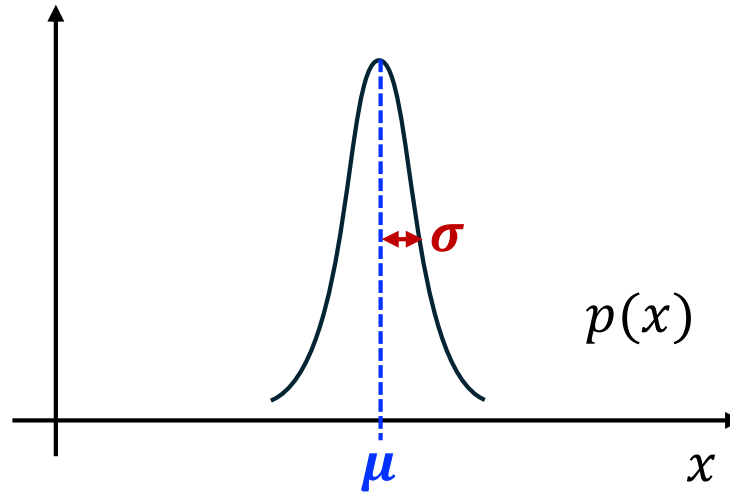


Gaussian distribution example

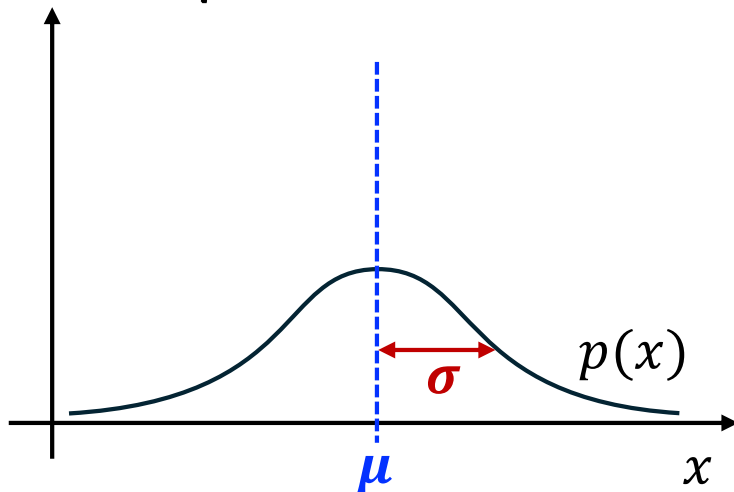
$$\mu = 0, \sigma = 1$$



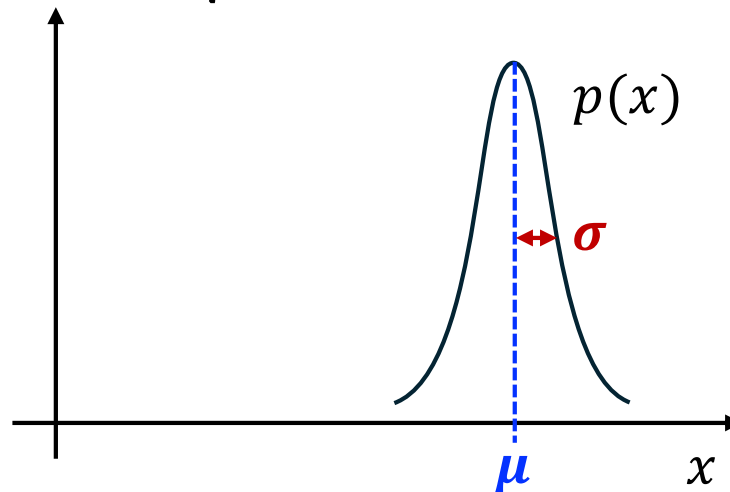
$$\mu = 0, \sigma = 0.5$$



$$\mu = 0, \sigma = 2$$

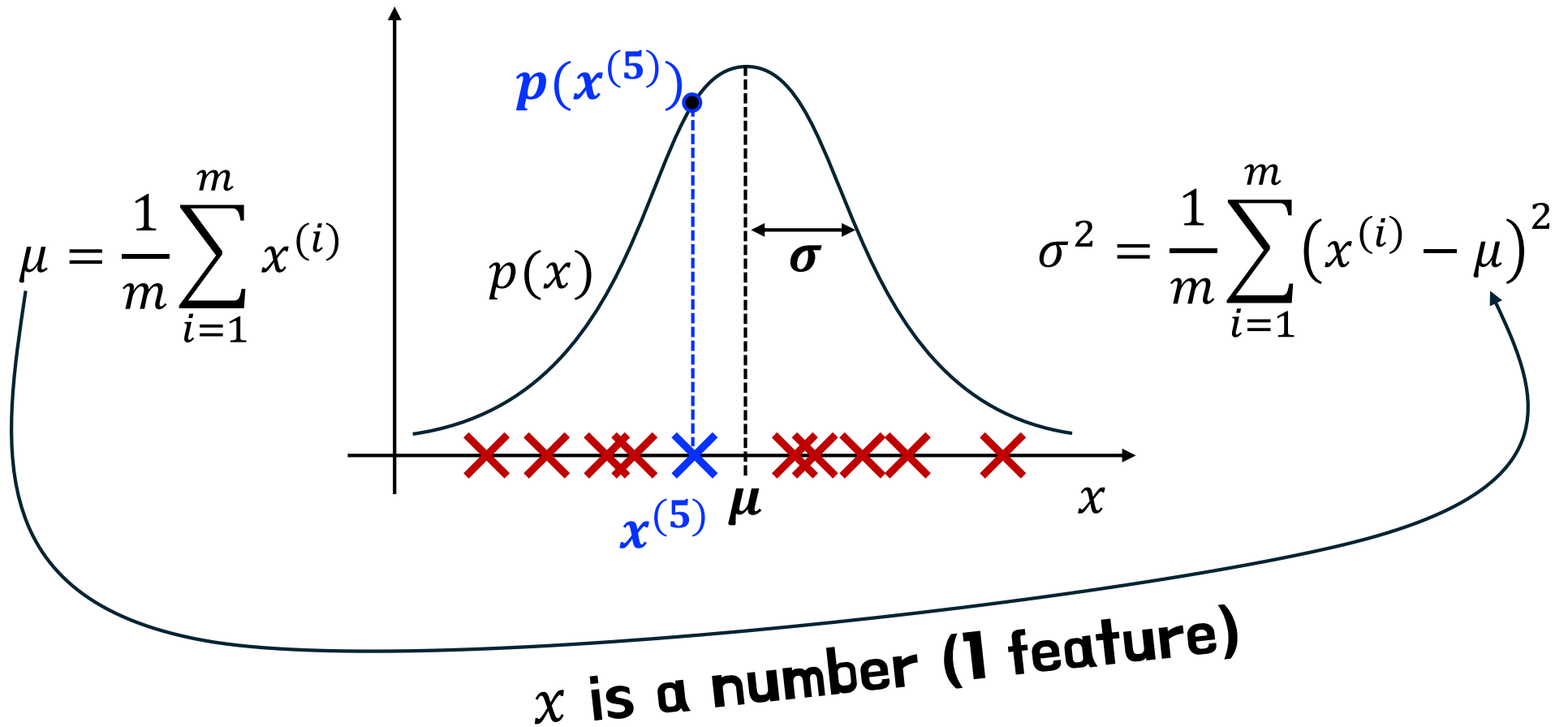


$$\mu = 3, \sigma = 0.5$$



Parameter estimation

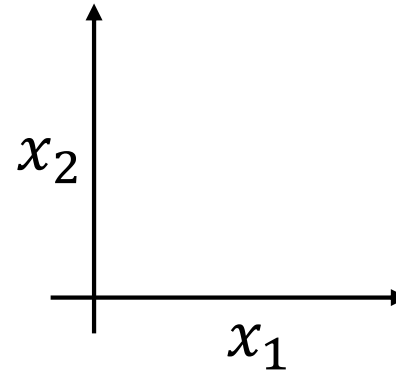
Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$



Density estimation

Training set: $\{\vec{x}^{(1)}, \vec{x}^{(2)}, \dots, \vec{x}^{(m)}\}$

Each example $\vec{x}^{(i)}$ has n features



$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$p(\vec{x}) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) * p(x_3; \mu_3, \sigma_3^2) * \dots * p(x_n; \mu_n, \sigma_n^2)$$

$$= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

$$p(x_1 = \text{high temp}) = \mathbf{1/10}$$

$$p(x_2 = \text{high vib.}) = \mathbf{1/20}$$

$$p(\vec{x}) = p(x_1) * p(x_2) = \mathbf{1/200}$$

Anomaly detection algorithm

1. Choose n features x_i that you think might be indicative of anomalous examples

2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

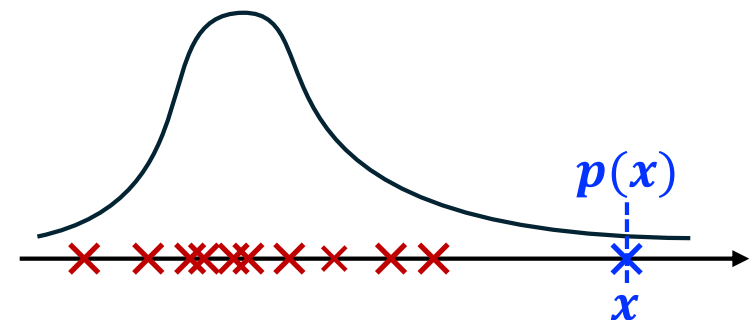
Vectorized formula

$$\vec{\mu} = \frac{1}{m} \sum_{i=1}^m \vec{x}^{(i)} \quad \vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

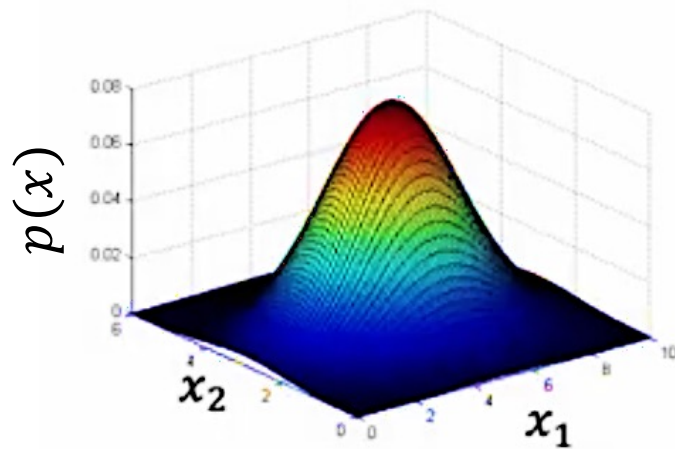
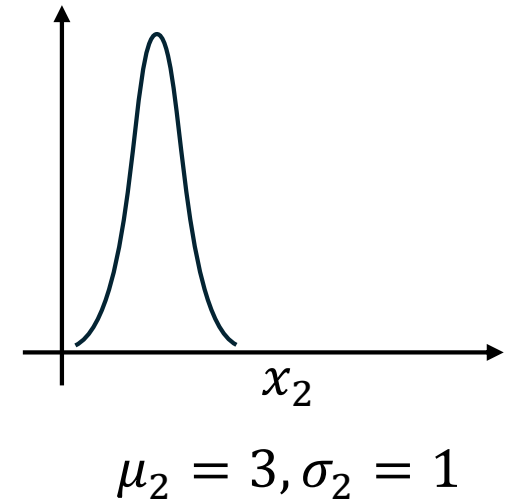
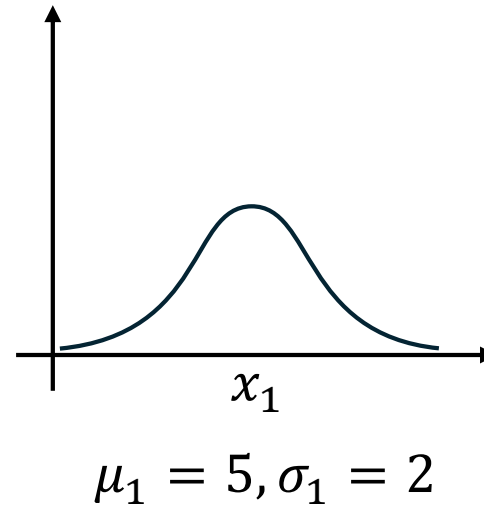
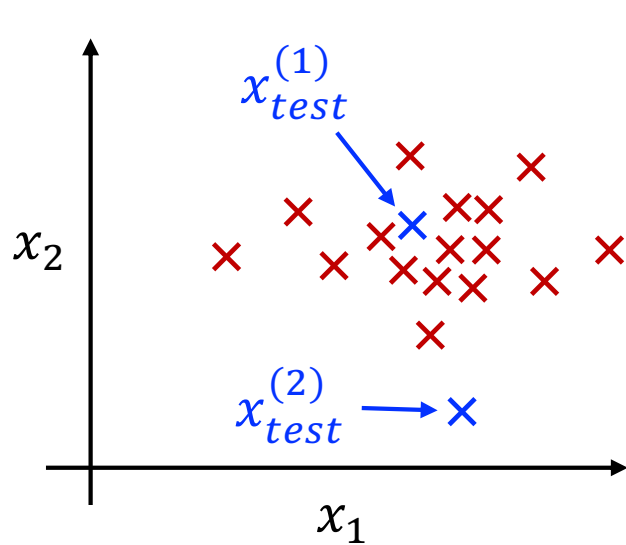
3. Given new example x , compute $p(x)$:

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}}$$

Anomaly if $p(x) < \varepsilon$



Anomaly detection example



$$\varepsilon = 0.02$$

$$p(x_{test}^{(1)}) = 0.0426$$

OK

$$p(x_{test}^{(2)}) = 0.0021$$

Anomaly



The importance of real-number evaluation

- When developing a learning algorithm (choosing features, etc.), making decision is much easier if we have **a way of evaluating our learning algorithm.**

- Assume we have some labeled data, of **anomalous** and **non-anomalous** examples.
 $y = 1$ $y = 0$

- **Training set:** $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume normal examples/not anomalous)
 $y = 0$ for all training examples

- **Cross validation set:** $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

- **Test set:** $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

Include a few anomalous examples

$(y = 1)$

Mostly normal examples

$(y = 0)$



Aircraft engines monitoring example

10,000 : Good (normal) engines

20 : Flawed engines (anomalous)

Training set : **6,000** good engines ($y = 0$) → **Train algorithm on training set**

CV : **2,000** good engines ($y = 0$), **10** anomalous ($y = 1$) → **Tune ε and x_j**

Test : **2,000** good engines ($y = 0$), **10** anomalous ($y = 1$)

If very few labeled anomalous examples,

Alternative: No test set

Training set : **6,000** good engines ($y = 0$)

CV : **4,000** good engines ($y = 0$), **2** anomalous ($y = 1$)

But, higher risk of overfitting



Algorithm evaluation

- Fit model $p(x)$ on training set $x^{(1)}, x^{(2)}, \dots, x^{(m)}$
- On a cross validation/test example x , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \varepsilon \text{ (normal)} \end{cases}$$

- Possible evaluation metrics:
 - True positive, false positive, false negative, true negative
 - Precision/Recall
 - F_1 -score
- Use cross validation set to choose parameter ε



Anomaly detection vs. supervised learning

Anomaly detection

- Very small number of positive examples ($y = 1$). **(0-20 is common)**
- Large number of negative examples ($y = 0$).
- Many different “types” of anomalies.
- Hard for any algorithm to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we've seen so far. **Fraud**

Supervised learning

- Large number of positive and negative examples. **(20 positive examples)**
- Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set. **Spam**



Anomaly detection vs. supervised learning

Anomaly detection

- Fraud detection
- Manufacturing - finding new previously unseen defects in manufacturing.
(e. g. aircraft engines)
- Monitoring machines in a data center

⋮

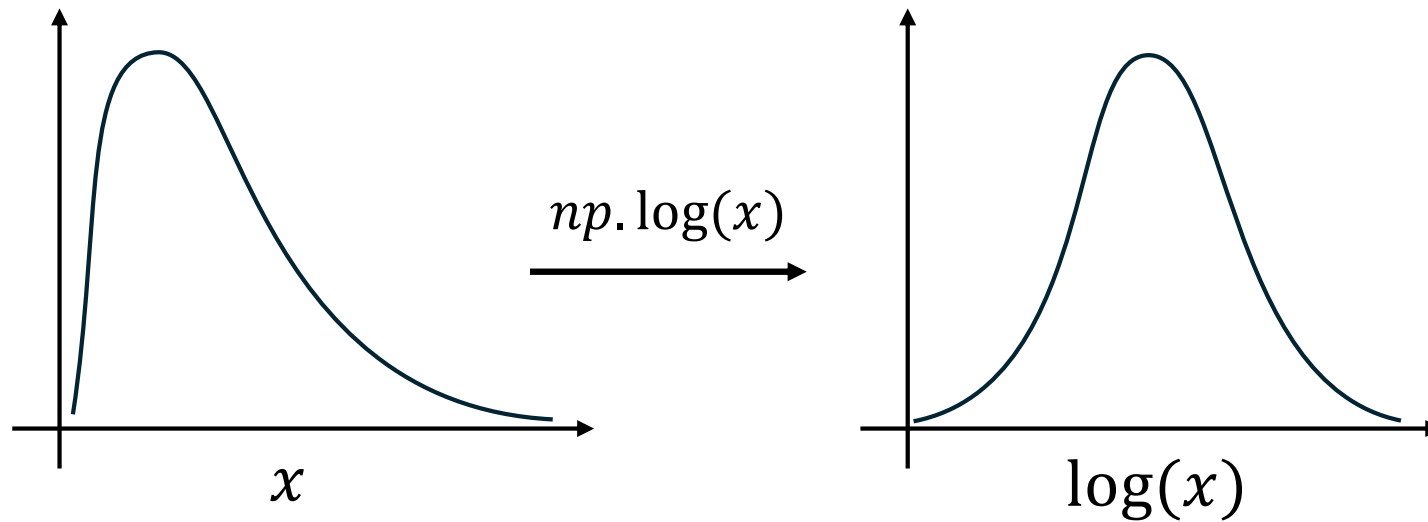
Supervised learning

- E-mail spam classification
- Manufacturing - finding known, previously seen defects
- Weather prediction (sunny/rainy/etc.)
- Diseases classification

⋮



Non-gaussian features



Other ways

- $x_1 \rightarrow \log(x_1)$
- $x_2 \rightarrow \log(x_2 + C)$
- $x_3 \rightarrow \sqrt{x_3}$
- $x_4 \rightarrow x_4^{1/3}$

Error analysis for anomaly detection

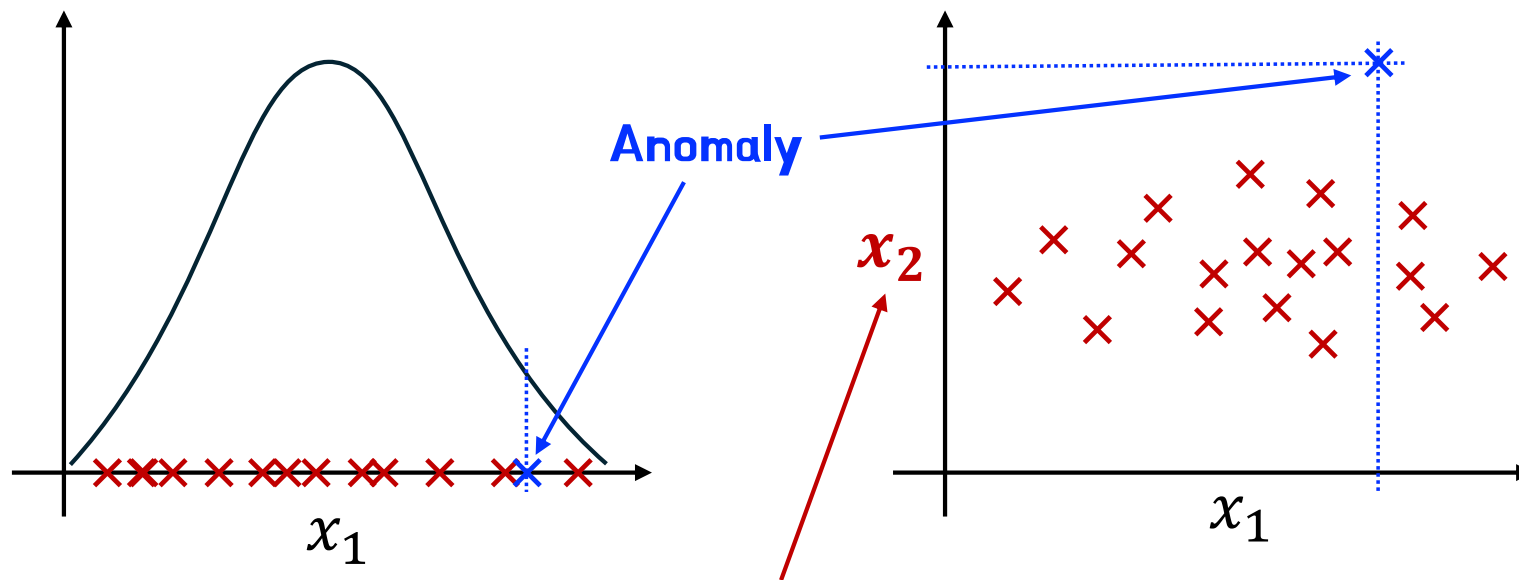
Want $p(x) > \varepsilon$: large for normal examples x .

$p(x) \leq \varepsilon$: small for anomalous examples x .

Most common problem:

$p(x)$ is comparable for normal and anomalous examples.

($p(x)$ is large for both)



Need to find the better feature

Monitoring computers in a data center

- Choose features that might take on unusually large or small values in the event of an anomaly.

$x_1 =$ memory use of computer

$x_2 =$ number of disk accesses/sec

$x_3 =$ CPU load

$x_4 =$ network traffic

$x_5 =$ CPU load / network traffic

$x_6 =$ (CPU load)² / network traffic

- Deciding feature choice based on $p(x)$
 - Large for normal examples;
 - Becomes small for anomaly in the cross validation set



Summary

- **Clustering (Unsupervised Learning)**

 - Group unlabeled data into meaningful clusters based on similarity.

- **K-means Algorithm**

 - Iteratively assign points to nearest centroids and update centroids to minimize distortion (cost function).

- **Choosing K & Initialization**

 - Proper K selection depends on purpose; random initialization requires multiple runs.

- **Anomaly Detection**

 - Model normal data distribution and flag low-probability examples as anomalies.

- **Anomaly Detection vs. Supervised Learning**

 - Use anomaly detection when anomalies are rare or unseen; use supervised learning when sufficient labeled data exists.

