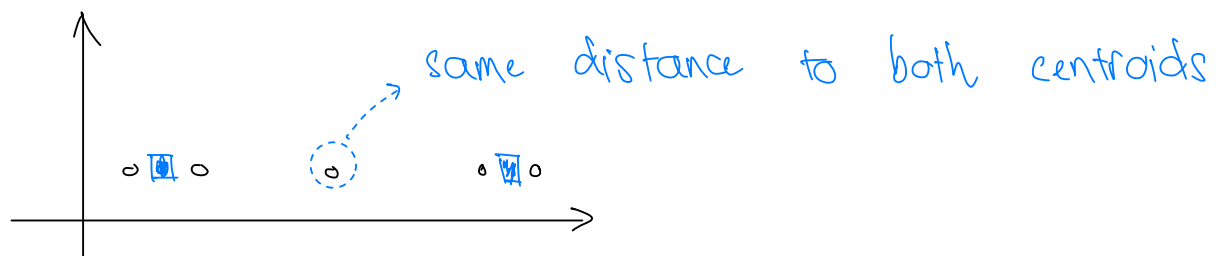# EXPECTATION MAXIMIZATION
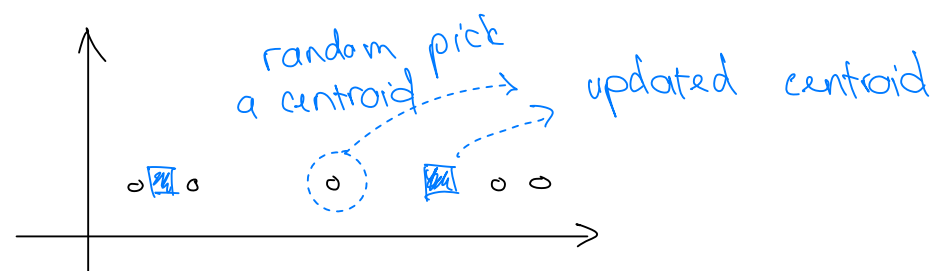
- Expectation Minimization (EM) a more generalized version of K-means
- Problem with K-means:
  - Iteration $t$ :

  same distance to both centroids

  - Iteration $t+1$ :

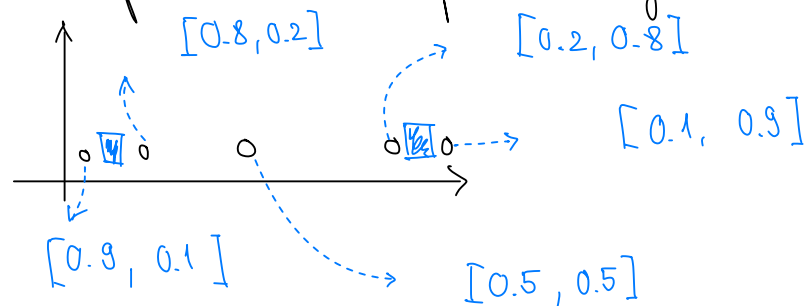  random pick a centroid → updated centroid

  $\Rightarrow$ Not ideal. The centroids are not symmetric even though the data is symmetric

- How EM solve that problem:

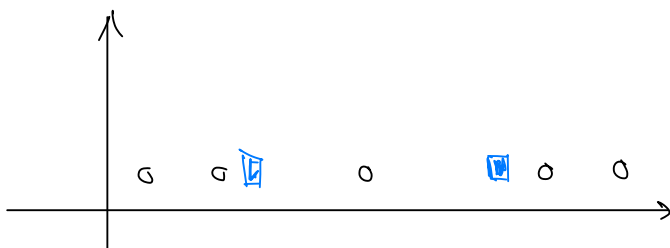  - In EM, each data point is a probability

  Iteration $t$:

  $[0.8, 0.2]$    $[0.2, 0.8]$    where:

  $[0.1, 0.9]$

  $[0.9, 0.1]$    $[0.5, 0.5]$

  $[0.9 \quad , \quad 0.1]$

  $P(z_i = 1)$       $P(z_i = 2)$

  - Since $z_i$ is a probability, $\mu$ will be more reprent the data

  $$\mu_1 = \frac{0.9 \, x_1 + 0.8 \, x_2 + 0.5 \, x_3 + 0.2 x_4 + 0.1 \, x_5}{0.9 + 0.8 + 0.5 + 0.2 + 0.1}$$

  $$\mu_2 = \frac{0.1 \, x_1 + 0.2 \, x_2 + 0.5 \, x_3 + 0.8 \, x_4 + 0.9 \, x_5}{0.1 + 0.2 + 0.5 + 0.8 + 0.9}$$

  $\Rightarrow$ Look much better! (more symmetrical)

# Expectation Maximization vs K-means

## K-means (Deterministic approach)

- **Assignment step:**

  For each data point $x_i$:
  $$z_i = \underset{k=1,\dots,K}{\arg\min} \| x_i - \mu_k \|^2$$

- **Centroids step:**

  For each centroid:
  $$\mu_k = \frac{\sum_{i \in S_k} x_i}{|S_k|} \quad \left( S_k = \{ i : z_i = k \} \right)$$

  $$= \frac{\sum^n I(z_i = k)\, x_i}{\sum^n I(z_i = k)}$$

## EM (Probabilistic approach)

- **E-step (Evaluation)**

  For each data point $x_i$:
  $$\begin{cases} P(z_i = 1) \\ \quad \vdots \\ P(z_i = K) \end{cases}$$
  $$\Rightarrow P(z_i = k) \quad \text{for each } k$$

  Example: sigmoid func distance
  $$P(z_i = k) = \frac{\exp\left(-\frac{1}{\lambda} \| x_i - \mu_k \|^2\right)}{\sum_{k=1}^{K} \left( -\frac{1}{\lambda} \| x_i - \mu_k \|^2 \right)}$$

- **M-step: (Maximization)**

  For each centroid:
  $$\mu_k = \frac{\sum^n P(z_i = k)\, x_i}{\sum^n P(z_i = k)}$$

---

## How to choose correct $\lambda$ ?

To find $\theta$, we use this procedure (Probabilistic Approach to Clustering)

### Probabilistic Approach to Clustering

- Assume some "hidden" joint distribution $p(x, z \mid \theta)$ that generates the data $x$ and the labels $z$.
- The goal is to find that distribution by estimating $\theta$
- To estimate $\theta$, there are 2 scenarios:

  i) $x$ and $z$ are known (complete information). Estimate $\theta$ by maximizing the <u>joint probability</u>
  $$\max_{\theta} \sum_{i=1}^{n} \log p(x_i, z_i \mid \theta)$$

  ii) Only $x$ are known (incomplete information). Estimate $\theta$ by maximizing the <u>marginal probability</u>
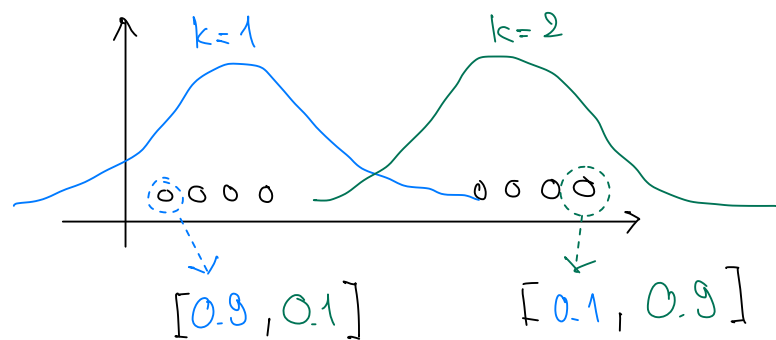  $$\max_{\theta} \sum_{i=1}^{n} \log p(x \mid \theta)$$

- Remember that the goal of clustering is to find the best possible assignment / label for each data point.
- Once we know the joint distribution $p(x, z \mid \theta)$ by ==estimating $\theta$== we can ==infer the label== of each data point.

$$z \sim p(z \mid x, \theta)$$

    observed →     → estimated earlier

## Example with Gaussian Distribution

- <u>Visual:</u>



$k=1$      $k=2$

$[0.9, 0.1]$      $[0.1, 0.9]$

- <u>Write down the joint distribution $p(x, z)$</u>

$$\underbrace{p(x, z)}_{\text{joint distribution}} = \underbrace{p(z)}_{\substack{\text{marginal} \\ \text{distribution}}} \cdot \underbrace{p(x \mid z)}_{\substack{\text{conditional} \\ \text{distribution}}}$$

For a specific $k$ and data point $x_i$:

$$\underbrace{p(x_i, z=k)}_{\substack{\text{joint probability data} \\ \text{and } z=k}} = p(z=k) \cdot p(x_i \mid z=k)$$

$$= \underbrace{\pi_k}_{\substack{\text{probability} \\ z = k}} \cdot \underbrace{N(x_i \mid \mu_k, \sigma_k^2)}_{\substack{\text{probability seeing } x \\ \text{given } z=k \text{ (likelihood)}}}$$

> **Trick:**
> Distribution is big if-else where each statement is a probability

==The goal== is to estimate $\theta$, so that we can ==infer the labels==

$$\theta = \{ \pi_k, \mu_k, \sigma_k^2 \}_{k=1}^{K}$$

*The nitty gritty of how to estimate $\theta$*

## Scenario 1:   Complete Information

○ We observed both data and labels (complete information)

○ Estimate $\theta$ by maximizing joint probability:
$$\hat{\theta} = \underset{\theta}{\arg\max} \sum_{i=1}^{n} \log p(x, z \mid \theta)$$

○ Detail calculations:

### i) Write down log-joint probability

For a ==specific data point and specific $k$==, the joint probability is:
$$p(x_i, z_i = k) = \pi_k \cdot N(x_i \mid \mu_k, \sigma_k^2)$$

Then for ==all data points and all values of $k$==:
$$\sum_{i=1}^{n} \sum_{k=1}^{K} I(z_i = k) \; p(x_i, z_i = k)$$
$$= \sum_{i=1}^{n} \sum_{k=1}^{K} I(z_i = k) \left[ \pi_k \cdot N(x_i \mid \mu_k, \sigma_k^2) \right]$$

Finally, arrive at log-joint probability:
$$\sum_{i=1}^{n} \sum_{k=1}^{K} I(z_i = k) \log \left[ \pi_k \cdot N(x_i \mid \mu_k, \sigma_k^2) \right]$$

### ii) Estimate $\theta$
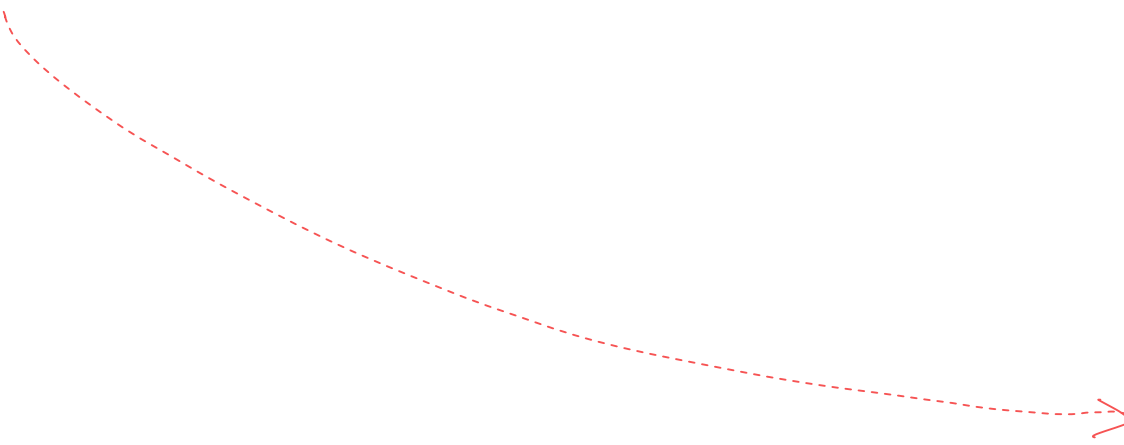
○ $\theta = \{ \pi_k, \mu_k, \sigma_k^2 \}_{k=1}^{K}$

○ Taking derivative w.r.t $\pi_k, \mu_k, \sigma_k$, we arrive at these results:

$$
\begin{cases}
\pi_k = \dfrac{\sum_{i=1}^{n} I(z_i = k)}{n} \\[3mm]
\mu_k = \dfrac{\sum_{i=1}^{n} I(z_i = k) \, x_i}{\sum_{i=1}^{n} I(z_i = k)} \\[3mm]
\sigma_k = \mathrm{Var}\left( \{ x_i \mid z_i = k \} \right)
\end{cases}
$$

→ variance of all data points where $z = k$

Obviously, in practice we won't always have labels, that is what we will analyze next

<u>Scenario 2:</u>   Incomplete information

- Can only observe data points, labels are hidden
- Estimating $\theta$ by maximizing marginal probability
$$\hat{\theta} = \underset{\theta}{\arg\max} \sum_{i=1}^{n} \log p(x_i \mid \theta)$$

- Detail calculations:

  i) <u>Write down log-marginal probability</u>

  For a specific data point:
  $$p(x_i \mid \theta) = \sum_{k=1}^{K} p(x_i, z_i = k \mid \theta)$$
  $$= \sum_{k=1}^{K} \pi_k \cdot N(x_i \mid \mu_k, \sigma_k^2)$$

  For all data points:
  $$\sum_{i=1}^{n} p(x_i \mid \theta) = \sum_{i=1}^{n} \sum_{k=1}^{K} \pi_k \cdot N(x_i \mid \mu_k, \sigma_k^2)$$

  Add the log:
  $$\sum_{i=1}^{n} \log p(x_i \mid \theta) = \sum_{i=1}^{n} \log\left(\sum_{k=1}^{K} \pi_k \cdot N(x_i \mid \mu_k, \sigma_k^2)\right)$$

  ii) <u>Estimate $\theta$:</u>

  - $\theta = \{\pi_k, \mu_k, \sigma_k^2\}$

  - There is no closed form solution to maximizing the log marginal probability

    Why? Lets analyze the marginal probability
    $$\sum_{i=1}^{n} \log\left(\sum_{k=1}^{K} \pi_k \cdot N(x_i \mid \mu_k, \sigma_k^2)\right)$$
    $$= \sum_{i=1}^{n} \log\left(\sum_{k=1}^{K} f(\pi_k, x_i, \mu_k, \sigma_k^2)\right)$$

    many distinct funtions
    $\iff$ many mountains

    - $\log(\text{moutain}) \iff$ steps to that mountain
    - $\log(\text{moutains}) \iff \Big[$ steps to mountain 1
      steps to mountain 2
      ....

  - So no closed form method, luckily we can still use Expectation Maximization (EM) to estimate $\theta$.

# Expectation Maximization Algorithm

○ <u>Pseudo code:</u>

- Initially $\theta$
- For $t = 1, 2, 3, \dots$

  i) <u>E - step (Evaluation):</u>

  Fill the hidden value $z_i$ by drawn i.i.d. from $p(z \mid x, \theta)$

  $$z_i \overset{i.i.d}{\sim} p(z \mid x, \theta)$$

  For specific data point $x_i$ and value $k$:

  $$p(z_i = k \mid x_i, \theta^t) = \frac{p(x_i, z_i = k \mid \theta^t)}{\sum_{k=1}^{K} p(x_i, z_i = k \mid \theta^t)}$$

  (posterior) conditional probability label of data point $x_i$ is $k$

  $\dashrightarrow$ prior · likelihood $= p(z_i = k) \cdot p(x_i \mid z_i = k)$

  marginal $p(x) = \sum^{K} p(x, z)$

  ii) <u>M - step (Maximization):</u>

  ○ Update $\theta$ by maximizing the expected log joint probability

  For specific data point and value $k$:

  $$\log p(x_i, z_i = k \mid \theta^t)$$

  For all data points and all values of $k$:

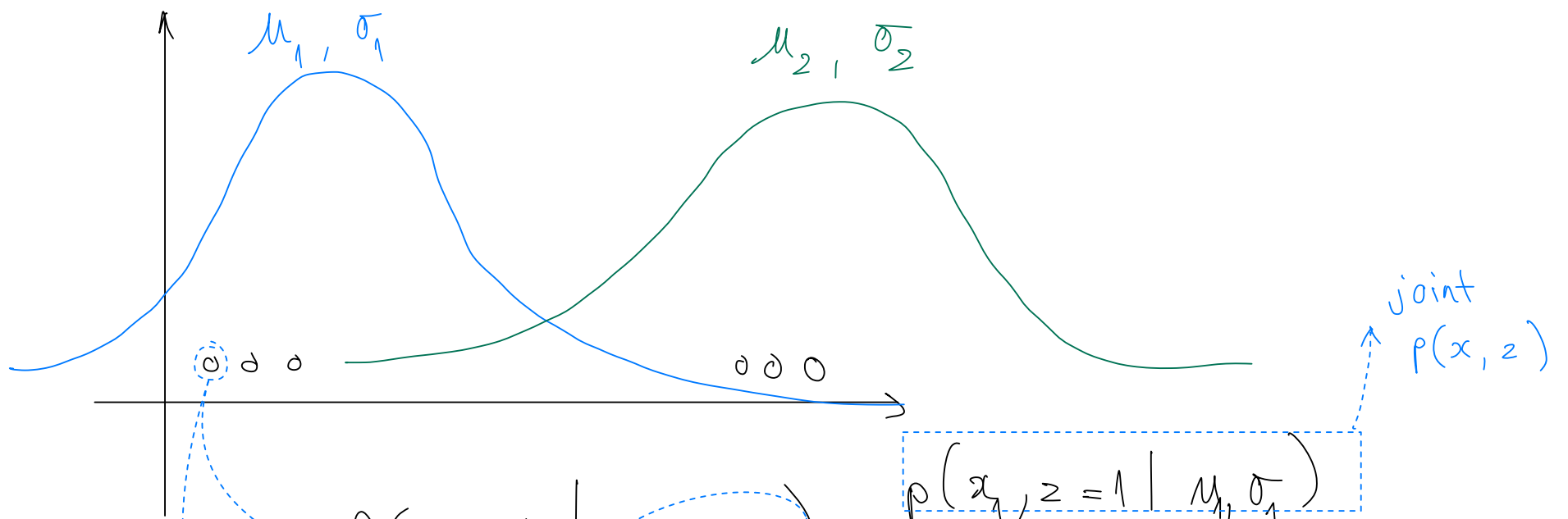  $$\sum_{i=1}^{n} \sum_{k=1}^{K} p(z_i = k \mid x_i, \theta^t) \log p(x_i, z_i = k \mid \theta^t)$$

  $I(z_i = k)$ in complete information case

  $$= \sum_{i=1}^{n} \sum_{k=1}^{K} p(z_i = k \mid x_i, \theta^t) \log \left[ \pi_k \cdot N(x_i \mid \mu_k, \sigma_k^2) \right]$$

  ○ Taking derivative w.r.t. $\pi_k, \mu_k, \sigma_k$, we arrive at these results:

  $$\begin{cases} \pi_k^{t+1} = \dfrac{\sum_{i=1}^{n} p(z_i = k \mid x_i, \theta^t)}{n} \\[3mm] \mu_k^{t+1} = \dfrac{\sum_{i=1}^{n} p(z_i = k \mid x_i, \theta^t) \cdot x_i}{\sum_{i=1}^{n} p(z_i = k \mid x_i, \theta^t)} \\[3mm] \sigma_k^{t+1} = Var(\{x_i \mid z_i = k\}) \end{cases}$$

  $\longrightarrow$ variance of data points with same label

$\mu_1, \sigma_1$   $\mu_2, \overline{\sigma_2}$

<u>Assignments</u>:

$$P\left(z = 1 \mid x_1, \mu_1, \sigma_1\right) = \frac{p\left(x_1, z = 1 \mid \mu_1, \sigma_1\right)}{\sum_{k}^{K} p\left(x_1, z = k \mid \mu_1, \sigma_1\right)}$$

joint $P(x, z)$

$\underbrace{\phantom{P(z=1\mid x_1,\mu_1,\sigma_1)}}_{\gamma_{11}}$    observed

marginal $P(x)$

$$P\left(z = 2 \mid x_1, \mu_2, \sigma_2\right)$$

$\underbrace{\phantom{P(z=2\mid x_1,\mu_2,\sigma_2)}}_{\gamma_{12}}$

Do the same for other data points, then:
we have the distribution of $z_i$

$$z_i \sim p\left(z \mid x, \mu, \sigma\right)$$

joint liklihood

<u>Centroids</u> $(\mu, \sigma)$:   $\{\mu, \sigma\}_{new} = \arg\max \sum \mathbf{E_{z_i \sim p(z \mid x, \mu, \sigma)}} \left[\log p\left(x_i, z_i \mid \mu_1, \sigma_1\right)\right]$

randomly drawn from