

Soft clustering:

- So far, clustering was done using **hard assignments** (1 point -> 1 cluster)
- Sometimes this doesn't accurately represent the data: it seems reasonable to have overlapping clusters.
- In this case, we can use **soft assignment** to assign points to every cluster with a certain probability.
  - Example:
    - Generate data using  $\mathbf{N}(\mu_1, \sigma_1)$  and  $\mathbf{N}(\mu_2, \sigma_2)$
    - we are given the weights of animals. Unknown to us these are weights from two different species.
    - Can we determine the species (group / assignment) from the height?
      - Any of these points could technically have been generated from either curve
      - We can create soft assignments based on these probabilities
- Mixture model:
  - X comes from a mixture model with k mixture components if the probability distribution of X is:

$$P(X = x) = \sum_{j=1}^k P(C_j)P(X = x|C_j)$$

- Gaussian mixture model:

$$\begin{aligned} \circ P(X = x|C_i) &\sim N(\mu, \sigma) \\ \circ P(X = x|C_i) &\sim N(\mu, \sigma) \end{aligned}$$

- GMM Clustering:

- **Goal:** Find the GMM that maximizes the probability of seeing the data we have.
- The probability of seeing the data we saw is (assuming each data point was sampled independently) the product of the probabilities of observing each data point.
- Finding the GMM means finding the parameters that uniquely characterize it. What are these parameters?
  - $\mathbf{P}(C_i)$  &  $\mu_i$  &  $\sigma_i$  for all  $k$  components.
  - Lets call  $\Theta = \{\mu_1, \dots, \mu_k, \sigma_1, \dots, \sigma_k, P(C_1), \dots, P(C_k)\}$
- Goal:

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^n \sum_{j=1}^k P(C_j)P(X_i | C_j)$$

- - Where  $\Theta = \{\mu_1, \dots, \mu_k, \sigma_1, \dots, \sigma_k, P(C_1), \dots, P(C_k)\}$
  - Joint probability distribution of our data
  - Assuming our data are independent
- How do we find the critical points of this function?

- Notice: taking the log-transform does not change the critical points
- Define :

$$l(\theta) = \log(L(\theta))$$

$$= \sum_{i=1}^n \log\left(\sum_{j=1}^k P(C_j)P(X_i | C_j)\right)$$

Expectation Maximization Algorithm:

- 1. Start with random  $\theta$
- 2. Compute  $\mathbf{P}(\mathbf{C}_j | \mathbf{X}_i)$  for all  $\mathbf{X}_i$  by using  $\theta$
- 3. Compute / Update  $\theta$  from  $\mathbf{P}(\mathbf{C}_j | \mathbf{X}_i)$
- 4. Repeat 2 & 3 until convergence