

Soft Clustering

1. Problem Statement

- Given a dataset of weights sampled from N different animals
- Can determine which weight belongs to which animal?

2. Output

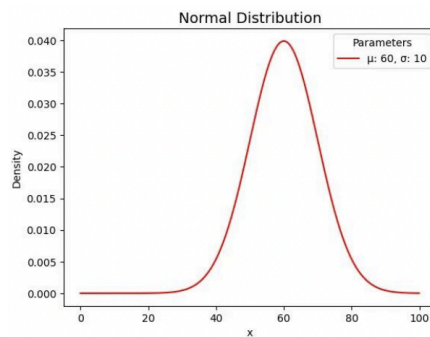
- Makes more sense to provide, for each data point (weight) the probability that it came from each species

$$P(S_j | X_i)$$

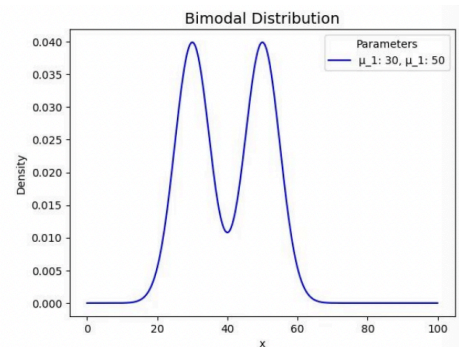
Where S_j is species j and X_i is the i th weight in the dataset

3. Things to Consider

- There is a prior probability of being one species (i.e. could have an imbalance dataset or there could just be more of one species than the other)
 - Ex: Some dinosaurs are more common than others: for example there are many Stegosauruses than Raptors in the park. This means a given data point, knowing nothing about it would just have a higher chance of being a Stegosaurus than a Raptor
- Weights vary differently depending on the species (i.e. each species could have a different weight distribution)



i.



4. How to Compute

$$P(S_j | X_i) = \frac{P(X_i | S_j)P(S_j)}{P(X_i)}$$

-
- $P(S_j)$ is the prior probability of seeing species S_j (that probability would be higher for Stegosauruses than the Raptors for example)
- $P(X_i | S_j)$ is the PDF of species S_j weights evaluated at weight X_i (seeing a Sauropod that weighs 100 tons is way more likely than seeing a Raptor that weighs 100 tons)

5. What about $P(X_i)$?

$$P(X_i) = \sum_j P(S_j)P(X_i|S_j)$$

a.

6. Mixture Model

a. X comes from a mixture model with k mixture components if the probability distribution of X is

$$P(X) = \sum_j P(S_j)P(X|S_j)$$

Mixture proportion
Represents the probability
of belonging to S_j

Probability of seeing x
when sampling from S_j

7. Gaussian Mixture Model

a. A Gaussian Mixture Model (GMM) is a mixture model where

$$P(X|S_j) \sim N(\mu, \sigma)$$

8. Maximum Likelihood Estimation (Intuition)

a. Suppose there is a given dataset of coin tosses and need to estimate the parameters that characterize that distribution - how would that be done?

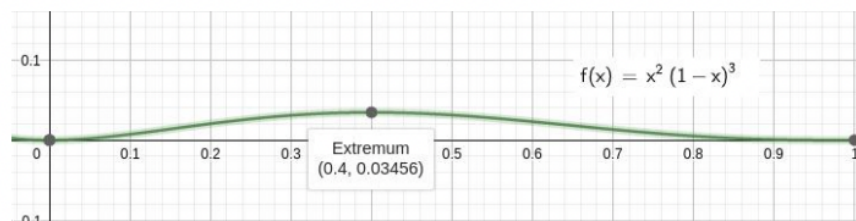
i. MLE: find the parameters that maximized the probability of having seen the data given

b. Example: Assume Bernoulli(p) iid coin tosses

Val
H
T
T
H
T

i. Find p that maximized that probability

ii. $P(\text{having seen the data we saw}) = p^2 * (1-p)^3$



iii.

iv. The sample proportion % is what maximizes this probability

9. GMM Clustering

- Goal: Find the GMM that maximizes the probability of seeing the data gathered

$$P(X_i) = \sum_j P(S_j)P(X_i|S_j)$$

- Recall:
- Finding the GMM means finding the parameters that uniquely characterize it.
- Parameters
 - $P(S_j)$ & μ_j & σ_j for all k components
 - Let's call $\Theta = \{\mu_1, \dots, \mu_k, \sigma_1, \dots, \sigma_k, P(S_1), \dots, P(S_k)\}$
- 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
- Goal:

$$\prod_i P(X_i) = \prod_i \sum_j P(S_j)P(X_i|S_j)$$

- How do we find the critical points of this function?
 - Take the log-transform since it does not change the critical points

$$\log \left(\prod_i \sum_j P(S_j)P(X_i|S_j) \right) = \sum_i \log \left(\sum_j P(S_j)P(X_i|S_j) \right)$$

ii.

$$\hat{\mu}_j = \frac{\sum_i P(S_j|X_i)X_i}{\sum_i P(S_j|X_i)}$$

$$\hat{\Sigma}_j = \frac{\sum_i P(S_j|X_i)(X_i - \hat{\mu}_j)^T(X_i - \hat{\mu}_j)}{\sum_i P(S_j|X_i)}$$

$$\hat{P}(S_j) = \frac{1}{N} \sum_i P(S_j|X_i)$$

iii.

10. Expectation Maximization Algorithm

- Start with random $\mu, \Sigma, P(S_j)$
- Compute $P(S_j | X_i)$ for all X_i by using $\mu, \Sigma, P(S_j)$
- Compute / Update $\mu, \Sigma, P(S_j)$ from $P(S_j | X_i)$
- Repeat 2 & 3 until convergence