

---

## Discussion 4

---

**Note:** Your TA will probably not cover all the problems on this worksheet. The discussion worksheets are not designed to be finished within an hour. They are deliberately made slightly longer so they can serve as resources you can use to practice, reinforce, and build upon concepts discussed in lectures, discussions, and homework.

**This Week's Cool AI Demo/Video:**

<https://oasis.decart.ai/starting-point>

### 1 Gaussian Mixture Clustering with Kangaroos and Berkeley Students (1D, Two Gaussians)

Consider a Gaussian Mixture Model (GMM) with two components in one dimension. Let  $z_n \in \{K, B\}$  be the latent (unobserved) class of the  $n$ -th jumper ( $K = \text{Kangaroo}$  and  $B = \text{Berkeley student}$ ) and  $x_n \in \mathbb{R}$  be the observed jump height (in meters).

Assume that the classes occur with equal probability and that, conditional on the class, jump heights are Gaussian with class-dependent means  $\mu_K$  and  $\mu_B$  (the average jump heights of Kangaroos and Berkeley Students, respectively), and a shared variance  $\sigma^2$  (the variability in jump heights).

(a) Write down the following quantities explicitly.

- (i) The **prior** probability that a jumper is a Kangaroo or Berkeley student before observing the jump height:  $p(z_n)$ .
- (ii) The **likelihood** of the jump height conditioned on whether the jumper is a Kangaroo or Berkeley student:  $p(x_n | z_n)$ .
- (iii) The **marginal** probability that a certain jump height is observed:  $p(x_n)$ .
- (iv) The **posterior** probability of a jumper being a Kangaroo or Berkeley student after observing the jump height:  $p(z_n | x_n)$ .

(b) Recall the K-means objective:

$$J(\{r_{nk}\}, \{\mu_k\}) = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

where  $\mu_k$  is the mean of cluster  $k$  and  $r_{nk} \in \{0, 1\}$  is a binary indicator variable that describes which of the  $K$  clusters the data point  $x_n$  is assigned to. The K-means algorithm alternates between updating the cluster assignments and updating the cluster centers.

We modify this algorithm to generate GMMs in the following way:

- (1) **Update assignments:** Replace the hard assignments  $r_{nk}$  with soft assignments  $\gamma_{nk}$ , representing the probability that a point  $x_n$  belongs to cluster  $k$ .
- (2) **Update parameters:** Rather than selecting cluster centers that minimize the squared distances of data points to their assigned cluster centers, select the parameters  $\theta$  that maximize the  $Q$  function:

$$Q(\theta) = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \ln p(x_n, z_n = k \mid \theta).$$

- (i) Identify which probability from part (a) corresponds to  $\gamma_{nk}$ . Verify  $\sum_k \gamma_{nk} = 1$ .
- (ii) Show that the  $Q$  function is actually the expected log-likelihood of our observed data,  $X = (x_1, \dots, x_N)$ , and the (unobserved) latent vector,  $Z = (z_1, \dots, z_N)$  under the posterior probability of  $Z$ .<sup>1</sup>
- (iii) Derive the maximum likelihood estimate for  $\mu_k$  using the *expected* log-likelihood,  $Q(\theta)$ , treating  $\gamma_{nk}$  as constants.

---

<sup>1</sup>We optimize the joint likelihood because it is much more tractable to compute while still optimizing the marginal likelihood. See this page for more information.

## 2 Closed-Form Solution for $\ell_2$ -Regularized Least Squares (Ridge)

Consider a data matrix  $X \in \mathbb{R}^{N \times D+1}$  with rows  $x_n^\top$ , a target vector  $t \in \mathbb{R}^N$ , and a parameter vector  $w \in \mathbb{R}^{D+1}$ . For a regularization hyperparameter  $\lambda \geq 0$ , define the ridge objective:

$$E(w) = \|Xw - t\|_2^2 + \lambda \|w\|_2^2.$$

- (a) Write  $E(w)$  in expanded quadratic form  $w^\top Aw - 2b^\top w + c$  by identifying  $A$ ,  $b$ , and  $c$  in terms of  $X$ ,  $t$ , and  $\lambda$ .

- (b) Compute the gradient  $\nabla_w E(w)$  using the identities

$$\nabla_w \|Xw - t\|_2^2 = 2X^\top(Xw - t), \quad \nabla_w \|w\|_2^2 = 2w.$$

- (c) Set the gradient to zero and derive the *normal equations* for ridge regression. Solve for  $w$ .

- (d) Justify why  $X^\top X + \lambda I$  is invertible for  $\lambda > 0$  (and discuss when it might fail for  $\lambda = 0$ ).