

ECE 271A: Statistical Learning I Quiz Report

JunPyung Kim
PID: A69044427
University of California, San Diego

December 3, 2025

1 Quiz 6: Gaussian Mixture Models for Cheetah Segmentation

1.1 Objective

In this quiz we revisit the cheetah image segmentation task using Gaussian mixture models (GMMs) estimated by the EM algorithm. The goals are:

- to study how sensitive the EM-trained GMMs are to random initialization when the number of components is fixed at $C = 8$, and
- to analyze how the number of mixture components $C \in \{1, 2, 4, 8, 16, 32\}$ affects the probability of error as we vary the feature dimension $d \in \{1, 2, 4, 8, 16, 24, 32, 40, 48, 56, 64\}$.

1.2 Methodology

1.2.1 Data and Feature Extraction

The training data are the DCT features provided in `TrainingSamplesDCT_8_new.mat`. There are 250 foreground (FG) samples and 1053 background (BG) samples, so the empirical class priors are

$$P(Y = \text{FG}) = \frac{250}{250 + 1053} \approx 0.192, \quad P(Y = \text{BG}) = \frac{1053}{250 + 1053} \approx 0.808.$$

The test image is the cheetah image `cheetah.bmp`. It is converted to grayscale and processed with a sliding 8×8 window (stride 1) over the entire 255×270 image. For each block we compute the 8×8 2-D DCT and then reorder the 64 coefficients using the zig-zag pattern given in `Zig-Zag Pattern.txt`. This produces 68850 feature vectors $\mathbf{x} \in \mathbb{R}^{64}$, one for each pixel position. The ground-truth mask is obtained from `cheetah_mask.bmp` by thresholding at 0.5.

1.2.2 Gaussian Mixture Model

For each class $Y \in \{\text{FG}, \text{BG}\}$ we model the conditional density of the full 64-dimensional feature vector as a mixture of C Gaussians with diagonal covariance matrices:

$$p(\mathbf{x} \mid Y = i) = \sum_{c=1}^C \pi_{i,c} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{i,c}, \Sigma_{i,c}), \quad (1)$$

where $\sum_c \pi_{i,c} = 1$ and $\Sigma_{i,c} = \text{diag}(\sigma_{i,c,1}^2, \dots, \sigma_{i,c,64}^2)$.

The models are trained with the EM algorithm. For a given class and C :

- **Initialization:** The means $\mu_{i,c}$ are initialized by selecting random training samples and adding small noise. All mixture weights are set to $\pi_{i,c} = 1/C$, and the initial diagonal variances are copied from the global variance of the training set (replicated across components).
- **E-step:** For each sample \mathbf{x}_n and component c we compute the log responsibility $\log \gamma_{n,c} = \log p(c \mid \mathbf{x}_n, Y = i)$ using the current parameters. This uses the diagonal Gaussian log-likelihood together with the mixture weights.
- **M-step:** From the responsibilities $\gamma_{n,c}$ we update

$$N_c = \sum_n \gamma_{n,c}, \quad \pi_{i,c} = \frac{N_c}{N},$$

$$\mu_{i,c} = \frac{1}{N_c} \sum_n \gamma_{n,c} \mathbf{x}_n, \quad \sigma_{i,c,d}^2 = \frac{1}{N_c} \sum_n \gamma_{n,c} (x_{n,d} - \mu_{i,c,d})^2 + \varepsilon,$$

where a small ε ensures numerical stability.

Iterations stop when the log-likelihood improvement becomes smaller than a tolerance threshold, or when the maximum number of iterations is reached.

1.2.3 Bayes Decision Rule and Error Computation

Once the FG and BG mixtures have been trained, we classify each image block using the Bayes decision rule

$$g(\mathbf{x}) = \log p(\mathbf{x} \mid Y = \text{FG}) + \log P(Y = \text{FG}) - \log p(\mathbf{x} \mid Y = \text{BG}) - \log P(Y = \text{BG}).$$

The pixel is labeled as cheetah if $g(\mathbf{x}) > 0$ and as grass otherwise. To study the effect of the feature dimensionality d , we always train the mixtures in the full 64-dimensional space, but at test time we keep only the first d zig-zag coefficients. This is implemented by restricting the mean and variance vectors to their first d entries when computing the log-likelihood.

Let \hat{Y} denote the classifier output and Y the true class (from the mask). The probability of error is estimated as

$$P_e = P(\hat{Y} \neq Y) = P(\hat{Y} \neq Y \mid Y = \text{FG})P(Y = \text{FG}) + P(\hat{Y} \neq Y \mid Y = \text{BG})P(Y = \text{BG}).$$

The conditional error terms correspond to the proportion of misclassified foreground and background pixels in the ground-truth mask.

1.3 Results and Discussion

1.3.1 Part (a): Effect of Initialization for $C = 8$ Components

For Part (a) we fix $C = 8$ components for both classes and train 5 independent mixtures for FG and 5 for BG using different random initializations. Combining them yields 25 classifier pairs; for each pair we compute the probability of error for all dimensions d in the set $\{1, 2, 4, 8, 16, 24, 32, 40, 48, 56, 64\}$. The resulting curves are shown in Figure 1.

Several qualitative trends can be observed:

- For all initialization pairs the error is highest when only the first DCT coefficient is used ($d = 1$) and remains relatively large for very low dimensions ($d \leq 8$), since a single or a few coefficients cannot capture the structure of the two classes.

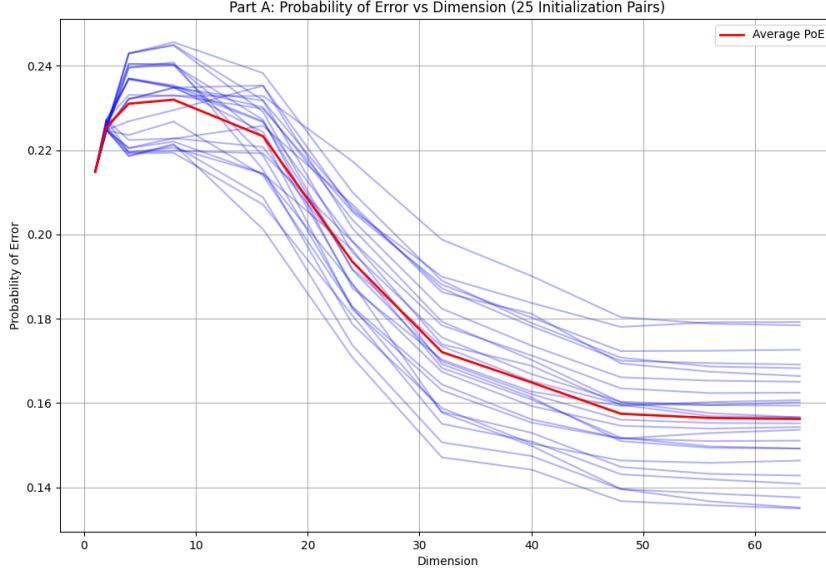


Figure 1: Part (a): probability of error vs. feature dimension for the 25 pairs of FG/BG mixtures obtained with different random initializations (blue curves) and their average (red curve).

- As more coefficients are included (d in the range roughly 16–40), the error drops significantly. The average curve shows a clear improvement when going from a handful of coefficients to a moderate-dimensional descriptor, reflecting that the mid-frequency DCT features are important for discriminating cheetah texture from grass.
- Beyond some dimension, gains become smaller and the curves tend to flatten. Some initializations even exhibit a slight increase in error at the largest dimensions, consistent with adding noisy or less informative coefficients and increasing estimation variance.
- The spread between the 25 curves is noticeable but not extreme: all runs follow the same global shape and have similar minima. This indicates that EM is somewhat sensitive to initialization (local maxima of the likelihood) but that the overall classifier performance is relatively robust: different initializations rarely change the error by more than a few percentage points.

Overall, Part (a) demonstrates that, although GMM training with EM does depend on the random start, the dominant factor affecting the probability of error is the feature dimension d : using only a very small number of coefficients is clearly sub-optimal, while using a moderate number of DCT features significantly improves segmentation quality.

1.3.2 Part (b): Effect of the Number of Components

For Part (b) we fix the initialization strategy and vary the number of mixture components

$$C \in \{1, 2, 4, 8, 16, 32\}$$

separately for FG and BG. For each C we train one GMM per class on the full 64-dimensional training data and then evaluate the probability of error as a function of the dimension d . The curves are shown in Figure 2.

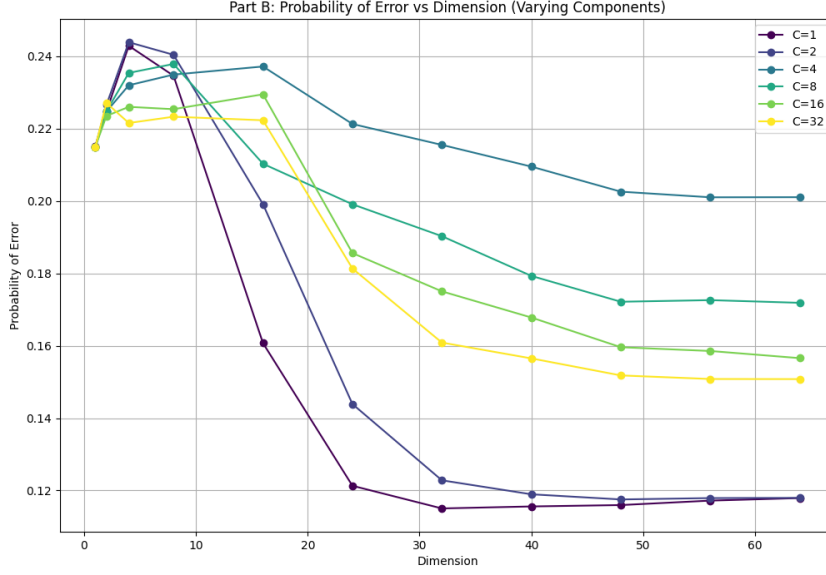


Figure 2: Part (b): probability of error vs. feature dimension for mixtures with different numbers of components $C \in \{1, 2, 4, 8, 16, 32\}$.

Some observations based on the values are:

- For very low dimensions ($d \leq 8$) all choices of C lead to very similar error rates (around 0.22–0.24). With so few features the representation is too limited for the additional mixture components to help.
- When more features are included, the models behave quite differently. For example, at $d = 32$ the estimated probability of error is approximately $P_e \approx 0.115$ for $C = 1$, 0.123 for $C = 2$, 0.216 for $C = 4$, 0.190 for $C = 8$, 0.175 for $C = 16$, and 0.161 for $C = 32$. The single-Gaussian model ($C = 1$) actually achieves the lowest error in this experiment.
- As C increases from 4 to 32 the curves tend to improve (e.g., the error for $C = 32$ at large d is lower than for $C = 4$ or $C = 8$), suggesting that additional components can help the model better fit the complex foreground and background distributions. However, none of the larger mixtures shows better performance than the simple $C = 1$ model on this task.
- For each fixed C , the error generally decreases as d grows from 16 to around 48–64, then slowly stabilizes. This is most evident for $C = 32$, whose error drops from about 0.22 at $d = 16$ to about 0.15 at $d = 56$ –64.

These results highlight that in theory increasing the number of mixture components should never hurt and should allow a more flexible approximation of the true class-conditional densities. In practice, with limited training data and diagonal covariances, the EM algorithm can overfit and converge to poor local optima when C is large. In our experiment the single-Gaussian model is strong, and adding more components often increases the probability of error instead of reducing it.

1.4 Conclusion

In this quiz we built GMM-based classifiers for cheetah vs. grass segmentation using DCT features of 8×8 image blocks. Part (a) showed that, while EM initialization causes some variation in performance, the overall shape of the probability-of-error curves is consistent across runs: using only a few DCT coefficients leads to high error, and using a moderate-to-large subset of coefficients substantially improves segmentation.

Part (b) demonstrated that simply increasing the number of mixture components does not guarantee better performance. With the available training data and diagonal covariances, a single Gaussian per class already models the DCT feature distributions well, and larger mixtures may overfit or get trapped in suboptimal local maxima of the likelihood. This illustrates the trade-off between model complexity and robustness in generative classification and emphasizes the need to validate mixture complexity on a separate performance metric such as segmentation error.