

# MixGoP\_Equation\_and\_Implementation

好的，我們來比對論文中 MixGoP 的數學方式與程式碼中的實作，主要會著重在 `main.txt` 中與 GMM 相關的部分。

論文中的 MixGoP 數學方式：

根據論文 "Leveraging Allophony in Self-Supervised Speech Models for Atypical Pronunciation Assessment\_v2.pdf" 和簡報 "ppt\_outline..." 以及總結 "summary\_by\_perplexity..."，MixGoP 的核心數學概念如下：

1. 為每個音位建模一個高斯混合模型 (GMM)。這個 GMM 旨在捕捉該音位的不同音位變體。
2. 音位可能性 (Phoneme Likelihood)  $P_{\theta}(s|p)$  的計算方式：

$$P_{\theta}(s|p) = \sum_{c=1}^C \pi_{c|p} \mathcal{N}(\text{Enc}(s) | \mu_{c|p}, \Sigma_{c|p})$$

其中：

- $C$  是每個音位 GMM 的子分布（成分）數量。
  - $\pi_{c|p}$  是音位  $p$  的第  $c$  個子分布的混合係數（權重）， $\sum_{c=1}^C \pi_{c|p} = 1$ 。
  - $\mathcal{N}(\text{Enc}(s) | \mu_{c|p}, \Sigma_{c|p})$  是以  $\mu_{c|p}$  為均值向量， $\Sigma_{c|p}$  為共變異數矩陣的高斯分布，用於評估語音片段的特徵  $\text{Enc}(s)$  在第  $c$  個子分布下的機率密度。 $\text{Enc}(s)$  代表從自監督語音模型 (S3M) 提取的語音片段  $s$  的特徵。
  - $\theta$  代表 GMM 的所有參數，包括  $\pi_{c|p}, \mu_{c|p}, \Sigma_{c|p} \in [C], p \in V$ ，其中  $V$  是音位集合。
3. MixGoP 分數 (MixGoP Score)  $\text{MixGoP}_p(s)$  的計算方式：

$$\text{MixGoP}_p(s) = \log P_{\theta}(s|p)$$

這個分數代表了給定音位  $p$  的情況下，觀察到語音片段  $s$  的對數似然度，分數越低表示發音越不典型。

程式碼 `main.txt` 中的實作 (根據您的描述)：

根據您之前的說明，`main.txt` 包含了訓練 GMM 的函數 `_train_gm`，並且使用了 `scikit-learn` 的 `GaussianMixture` 類別。讓我們來分析 `_train_gm` 的可能實作如何對應到上述數學公式：

1. GMM 建模：

- 程式碼中會針對每個音位，創建一個 `GaussianMixture` 物件。

- `n_components` 參數在 `GaussianMixture` 的初始化中被設定為一個常數（例如 32），這直接對應到數學公式中的  $C$ ，即每個音位 GMM 的子分布數量。
- `covariance_type` 參數（例如 `'full'`，`'diag'`，`'spherical'`，`'tied'`）決定了共變異數矩陣  $\Sigma_{c|p}$  的形式。
- `n_init` 參數控制了 GMM 訓練的初始化次數，這有助於找到更好的模型參數，但並不直接對應到最終的數學公式。
- `init_params` 參數通常設定為 `'kmeans'`，表示使用 K-means 演算法初始化 GMM 的均值  $\mu_{c|p}$ 。論文也提到使用 k-means 進行初始化。
- `max_iter` 參數設定了 EM 演算法的最大迭代次數，用於優化 GMM 的參數  $\pi_{c|p}, \mu_{c|p}, \Sigma_{c|p}$ 。然而，論文提到由於特徵的高維度，EM 優化步驟可能很少或沒有，這意味著初始化的中心可能與最終的中心相似。

## 2. 訓練資料和特徵：

- `_train_gm` 函數接收語音特徵 `embs` 作為輸入，這些特徵是由 `extract_features.txt` 從 S3M 模型中提取的  $Enc(s)$ 。
- 程式碼中可能使用了 `n_sample` 參數來限制每個音位用於訓練 GMM 的最大樣本數量，這與論文中討論的 MixGoP 的樣本效率有關。訓練資料會根據音位標註進行分割。

## 3. MixGoP 分數計算 (推測)：

- 雖然 `main.txt` 中可能沒有一個明確名為 `calculate_mixgop_score` 的函數，但在評估階段，對於一個測試語音片段的特徵向量  $Enc(s)$  和其對應的音位  $p$ ，程式碼會使用已經訓練好的音位  $p$  的 GMM 模型 (`gm[p]`) 來計算該特徵向量在這個 GMM 下的對數似然度 ( $\log\text{-likelihood}$ )。
- `scikit-learn` 的 `GaussianMixture` 物件提供了 `score_samples(X)` 方法，該方法可以計算每個樣本在模型下每個子分布的加權對數似然度。要得到  $P_{\theta}(s|p)$  的對數，可能需要對這些加權對數似然度進行操作（例如使用 `logsumexp` 來近似對數和）。或者，`GaussianMixture` 也可能直接提供計算整體對數似然度的方法。
- 這個計算出的對數似然度就對應到論文中的  $MixGoP_p(s) = \log P_{\theta}(s|p)$ 。

## 總結比對：

- 子分布數量 ( $C$ ): 程式碼中的 `n_components` 參數直接對應到數學公式中的  $C$ ，控制了每個音位 GMM 的複雜度，以捕捉多個音位變體。
- 混合係數 ( $\pi_{c|p}$ )、均值 ( $\mu_{c|p}$ )、共變異數 ( $\Sigma_{c|p}$ ): 這些 GMM 的參數是透過 `GaussianMixture` 的 `fit()` 方法，基於每個音位的 S3M 特徵進行估計的。訓練後的 `GaussianMixture` 物件會儲存這些參數。
- 高斯分布 ( $\mathcal{N}$ ): `GaussianMixture` 內部使用了高斯分布來建模每個子分布的機率密度。
- 特徵 ( $Enc(s)$ ): 程式碼接收由 S3M 模型提取的特徵作為訓練和評估 GMM 的輸入，這與數學公

式中的  $Enc(s)$  相符.

- **對數似然度 ( $\log P_\theta(s|p)$ ):** 程式碼在評估階段會使用訓練好的 GMM 模型來計算測試語音片段在對應音位模型下的對數似然度，這就是 MixGoP 的核心分數.

因此，`main.txt` 中的程式碼實作與論文中 MixGoP 的數學原理是高度一致的。程式碼使用 `scikit-learn` 提供的 `GaussianMixture` 類別來建立和訓練每個音位的 GMM，並利用 S3M 特徵作為輸入。在評估階段，通過計算測試語音在相應 GMM 下的對數似然度來獲得 MixGoP 分數，從而評估發音的非典型程度.