# Updating GMMs

Riddhanya, Tarun, Mitushi, Arun

November 2024

If the new Gaussian Mixture Model (GMM) is **unknown** and you only receive samples from this new distribution, you need to adjust your current GMM iteratively as new data arrives, without directly knowing the parameters of the new GMM. In this scenario, you can update your existing GMM to align with the new samples using a method like **online/incremental learning** or **adaptive GMM updating**.

Here's how you can handle this:

# 1 **Online (Incremental) GMM Updates**

In this method, you update the parameters of your current GMM incrementally as each new sample from the unknown GMM arrives. This approach adjusts the means, covariances, and weights of your GMM to account for the new data without retraining the entire model from scratch.

Here's an outline of how to perform incremental updates:

- **Assign the new sample to a component** of the GMM based on the posterior probability that the sample belongs to each component. - **Update the parameters** (mean, covariance, and weight) of the GMM component based on the new sample.

The update for each parameter can be done using a **learning rate** to control how quickly the parameters are updated.

**Example of Parameter Updates**   Given a sample $x_{\text{new}}$, if it is assigned to component $k$ of your GMM, you can update the parameters for that component $(\mu_k, \Sigma_k, w_k)$:

- **Mean update** for component $k$:

$$\mu_k^{\text{new}} = \mu_k^{\text{old}} + \alpha \cdot (x_{\text{new}} - \mu_k^{\text{old}})$$

Where $\alpha$ is the learning rate (typically small, such as 0.01).

- **Covariance update** for component $k$:

$$\Sigma_k^{\text{new}} = \Sigma_k^{\text{old}} + \alpha \cdot \left[ (x_{\text{new}} - \mu_k^{\text{old}})(x_{\text{new}} - \mu_k^{\text{old}})^T - \Sigma_k^{\text{old}} \right]$$

- **Weight update**:

The weight $w_k$ for component $k$ can be updated to reflect the fact that more data is now being assigned to that component:

$$w_k^{\text{new}} = w_k^{\text{old}} + \alpha \cdot (1 - w_k^{\text{old}})$$

The weights should still sum to 1, so after updating the weights, normalize them across all components.

This process can be repeated for each new sample you receive.

# 2 **Expectation-Maximization (EM) with Streaming Data**

If you want a more principled approach, you can apply the **Expectation-Maximization (EM) algorithm** in an online or batch mode where you re-estimate the GMM parameters in mini-batches as new data arrives.

**How to Implement Online EM:** - **E-step**: For each new mini-batch of data, calculate the responsibility (posterior probability) that each data point belongs to each GMM component. - **M-step**: Update the GMM parameters (means, covariances, and weights) based on the responsibilities calculated in the E-step.

To make it work with streaming data, you can run the EM algorithm over small batches of data rather than over the entire dataset. Each time a new batch of data is available, update the GMM with a **weighted combination** of the previous model and the new estimates.

—

# 3 **Adaptive GMM with New Component Creation**

If the new data differs significantly from your current GMM, you might need to **adapt** your GMM by adding new components. This method gradually grows the GMM as you observe more data.

**Steps:**

1. **Check model fit**: For each new sample or batch, calculate the likelihood of the sample under the current GMM. If the likelihood is too low (indicating that the current GMM doesn't fit the new data well), consider adding a new component.

2. **Add new components**: When a new component is needed, initialize its parameters based on the new data and gradually incorporate it into the GMM.

3. **Prune components**: To avoid overfitting, you can also prune components with very small weights after adding new components.

This method allows the model to adapt flexibly to new data without knowing the true parameters of the new distribution.

—

# 4 **Dirichlet Process Gaussian Mixture Models (DP-GMM)**

If you are uncertain about how many components the new distribution might have, you can use a **Dirichlet Process Gaussian Mixture Model (DP-GMM)**. This is a nonparametric Bayesian approach where the number of mixture components is not fixed and can grow dynamically as more data is observed.

- In a DP-GMM, new components are added automatically as needed based on the data, and the model adjusts itself to better represent the underlying distribution. - The **Dirichlet Process** acts as a prior over the number of components, allowing the number to grow as new data indicates the need for additional components.

This approach works well if you suspect that the new data might introduce new Gaussian components that were not present in the prior GMM.

—

**Example of Incremental Updates (Python Code):**

```python
import numpy as np
from sklearn.mixture import GaussianMixture

# Existing GMM (assumed)
gmm_prior = GaussianMixture(n_components=2)
prior_data = np.random.normal(loc=0, scale=1, size=(100, 2))
gmm_prior.fit(prior_data)

# New data stream (from unknown GMM)
new_samples = np.random.normal(loc=3, scale=2, size=(10, 2))  # Incoming samples

# Incremental update of GMM parameters
alpha = 0.05  # Learning rate for updates
```

```
for sample in new_samples:
    # Calculate the posterior probabilities (responsibilities) for the new sample
    responsibilities = gmm_prior.predict_proba(sample.reshape(1, -1))

    # Update the parameters of the GMM for each component
    for k in range(gmm_prior.n_components):
        # Update the mean for component k
        gmm_prior.means_[k] = gmm_prior.means_[k] + alpha * responsibilities[0, k] * (sample - gmm_prior.m

        # Update the covariance for component k (using a simplified diagonal covariance model)
        diff = sample - gmm_prior.means_[k]
        gmm_prior.covariances_[k] = gmm_prior.covariances_[k] + alpha * responsibilities[0, k] * (np.outer

    # Update the weights (normalizing afterwards)
    gmm_prior.weights_ = gmm_prior.weights_ + alpha * (responsibilities[0] - gmm_prior.weights_)
    gmm_prior.weights_ /= np.sum(gmm_prior.weights_)  # Ensure weights sum to 1

# The GMM model is now incrementally updated with the new samples
‘‘‘
```

# 5 Summary of Approaches:

1. **Incremental GMM Updates**: Gradually adjust the parameters of your GMM with each new sample. 2. **Online EM**: Use mini-batches of new data to iteratively update your GMM using the EM algorithm. 3. **Adaptive GMM**: Add new components as new data indicates the need for more complexity. 4. **Dirichlet Process GMM (DP-GMM)**: Let the model automatically decide how many components to use, growing as needed.

Let me know if you'd like more details on one of these methods or a specific code example!

# 6 Dirichlet Process Gaussian Mixture Model (DP-GMM)

A **Dirichlet Process Gaussian Mixture Model (DP-GMM)** is a nonparametric extension of a Gaussian Mixture Model (GMM) that allows for a **flexible number of components**. Unlike a traditional GMM where you specify a fixed number of components $K$, the DP-GMM adapts the number of components based on the data, potentially introducing new components as more data is observed.

**Key Concepts Behind DP-GMM:**   - **Dirichlet Process (DP)**: The Dirichlet Process is a distribution over distributions. It allows for a flexible number of mixture components, where the number of components grows as more data is added. - **Chinese Restaurant Process (CRP)**: The DP-GMM can be understood through the CRP metaphor, where each new data point either joins an existing cluster (component) or starts a new cluster, depending on a probability that balances the size of the existing clusters and a prior on starting new ones. - **Nonparametric Nature**: Instead of predefining the number of clusters $K$, the model can dynamically create and adjust the number of components as needed.

In DP-GMM, the number of components is theoretically infinite, but in practice, the model uses only as many components as are supported by the data.

—

**Mathematical Formulation:**   Given data $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_N$, the goal is to model the distribution of the data as a mixture of Gaussians, where the number of components is determined dynamically.

1. **Generative Model**: - Let $\pi \sim \text{StickBreakingProcess}(\alpha)$, where $\alpha$ is a concentration parameter. - For each data point $\mathbf{x}_i$, assign it to a component $z_i \sim \pi$. - The $k$-th component $z_i = k$ is associated with a Gaussian distribution $\mathcal{N}(\mu_k, \Sigma_k)$. - The concentration parameter $\alpha$ controls how likely new components are created. Larger $\alpha$ leads to more components.

2. **Inference**: - The parameters $\theta_k = (\mu_k, \Sigma_k)$ of each Gaussian component are learned from the data. - As new data arrives, the DP-GMM can adjust the number of components and reassign data points to existing or new components.

**Components:** - **Concentration Parameter $\alpha$**: Controls the tendency to add new clusters. Higher $\alpha$ encourages more clusters. - **Stick-Breaking Construction**: A common method to sample from a DP, where each new component is assigned a proportion of the remaining "stick" length, ensuring the sum of all weights is 1.

—

**Dirichlet Process Gaussian Mixture Model (DP-GMM) in Practice:** The DP-GMM can be implemented using libraries like 'scikit-learn', 'PyMC3', or 'numpyro'. Below is an example using 'scikit-learn', which provides an approximate implementation of the DP-GMM using the **Variational Inference (VI)** algorithm.

**DP-GMM Example using 'scikit-learn':**

```python
from sklearn.mixture import BayesianGaussianMixture
import numpy as np
import matplotlib.pyplot as plt

# Generate synthetic data from an unknown mixture of Gaussians
np.random.seed(42)
n_samples = 300

# Data is generated from 3 different Gaussian distributions
X1 = np.random.normal(loc=0, scale=1, size=(n_samples//3, 2))
X2 = np.random.normal(loc=5, scale=2, size=(n_samples//3, 2))
X3 = np.random.normal(loc=-5, scale=0.5, size=(n_samples//3, 2))

X = np.vstack((X1, X2, X3))

# DP-GMM implementation using scikit-learn
dp_gmm = BayesianGaussianMixture(
    n_components=10,  # Maximum number of components (can grow dynamically)
    covariance_type='full',  # Full covariance for each Gaussian
    weight_concentration_prior_type='dirichlet_process',  # Use Dirichlet Process prior
    weight_concentration_prior=0.1,  # Prior on the concentration parameter
    max_iter=1000,
    random_state=42
)

# Fit the DP-GMM model to the data
dp_gmm.fit(X)

# Predict the labels (component assignments) for each sample
labels = dp_gmm.predict(X)

# Plot the resulting clusters and their Gaussians
plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis', s=40, marker='o', edgecolors='k')

# Display the means of each Gaussian component
plt.scatter(dp_gmm.means_[:, 0], dp_gmm.means_[:, 1], s=200, c='red', marker='x')

plt.title("DP-GMM Clustering Result")
plt.show()
```

**Explanation of Key Parameters:** - 'n_components=10': The maximum number of components. DP-GMM will determine how many components are actually needed based on the data. - 'weight_concentration_prior_type='dirichlet_process'': This sets the prior to use the Dirichlet Process, which allows the model to grow components dynamically. - 'weight_concentration_prior=0.1': The concentration parameter controls how likely the model is to create new components. A smaller value encourages fewer components, while a larger value allows more components. - 'covariance_type='full'': Specifies that each Gaussian component should have a full covariance matrix.

**Key Points:** - **Automatic Component Adaptation**: The DP-GMM will automatically adjust the number of clusters (Gaussian components) as data is observed, without needing to predefine the number of components. - **Concentration Parameter $\alpha$**: Controls how likely the model is to create new components. A high $\alpha$ creates more clusters, while a low $\alpha$ encourages fewer clusters. - **Variational Inference (VI)**: 'scikit-learn' uses a variational method to approximate the DP-GMM, which is more efficient than directly sampling from the true posterior.

—

**Advantages of DP-GMM:** 1. **Dynamic Number of Components**: It adapts to the data, growing the number of Gaussian components as needed. 2. **Flexibility**: You don't need to specify the number of components beforehand, making it ideal for cases where the underlying distribution's complexity is unknown. 3. **Richer Clustering**: DP-GMM can capture more complex, multimodal data distributions.

**Disadvantages:** 1. **Complexity**: The nonparametric nature of the DP-GMM introduces additional complexity, both in terms of computation and interpretability. 2. **Approximation**: The variational inference used by 'scikit-learn' is an approximation of the true posterior, so results may not always match fully Bayesian approaches.

—

**When to Use DP-GMM:** - When you suspect the data may contain an unknown number of clusters or components. - When you want a model that can grow as more data is added. - For problems like density estimation, clustering, and anomaly detection where the underlying distribution is not well-known.

Would you like further assistance with any part of the implementation, or a deeper dive into any of the concepts?

# 7 Mathematical Guarantees

The Gaussian Mixture Model (GMM) update methods mentioned earlier—**Incremental Updates**, **Online Expectation-Maximization (EM)**, and **Dirichlet Process Gaussian Mixture Models (DP-GMM)**—each have varying levels of convergence guarantees, depending on the approach and the algorithm used. Here's a breakdown of the mathematical guarantees on convergence for each method:

## 7.1 **Incremental (Online) GMM Updates**

The method of updating GMM parameters incrementally using a learning rate (like in the example code I provided) is a heuristic, often employed in **streaming** or **nonstationary** environments. While this approach works well in practice for tracking changing distributions, it lacks rigorous convergence guarantees in the general case.

- **Convergence Guarantee**: Incremental updates based on a small learning rate might not have strong convergence guarantees, especially in nonstationary settings. The updates are sensitive to the learning rate, and too large of a rate may prevent convergence. - **In Practice**: For stationary data, if the learning rate $\alpha$ is chosen carefully (e.g., diminishing over time), the method can converge to a local optimum. However, it's not guaranteed to converge to the global optimum of the likelihood function.

## 7.2 **Online Expectation-Maximization (EM) Algorithm**

The **online version of the EM algorithm** is designed to handle streaming data and provides a principled way to update GMM parameters as new data comes in. There are some mathematical guarantees related to the convergence of this method:

| Method | Convergence Guarantee |
|---|---|
| Incremental GMM Updates | Heuristic method; no strong theoretical guarantees. May converge with a proper learning rate, but sensitive to tuning. |
| **Online EM** | Converges to a local optimum of the likelihood function under Robbins-Monro conditions on the learning rate. |
| **Standard EM** | Converges to a local maximum of the likelihood function, with monotonic likelihood improvement at each step. |
| **DP-GMM (Variational)** | Converges to a local optimum of the variational objective (ELBO); global convergence not guaranteed. |
| **DP-GMM (MCMC-based)** | Asymptotically samples from the true posterior; convergence rate can be slow, but theoretically converges to the true posterior. |

Table 1: Caption

**a. **Classical EM Convergence**:** - The **standard (batch) EM algorithm** for GMMs is guaranteed to converge to a **local optimum** of the likelihood function, but not necessarily to the global optimum. - EM monotonically increases the data likelihood (or leaves it unchanged) at each step, so it is guaranteed to converge to some point, though this point could be a local maximum or saddle point.

**b. **Online EM Convergence**:** - **Online EM** extends classical EM to handle streaming data using stochastic updates. - **Robbins-Monro conditions**: If the learning rate $\alpha_t$ follows the Robbins-Monro conditions (i.e., $\sum \alpha_t = \infty$ and $\sum \alpha_t^2 < \infty$), the online EM algorithm can converge almost surely to a stationary point of the likelihood function.

- **Convergence Guarantee**: Online EM has convergence guarantees to a local optimum under proper conditions on the learning rate and data distribution. Specifically, for the **stochastic approximation EM algorithm**, it can be shown to converge to a stationary point of the expected log-likelihood, but global convergence is still not guaranteed.

## 7.3 **Dirichlet Process Gaussian Mixture Model (DP-GMM)**

The **DP-GMM** uses nonparametric Bayesian inference to dynamically adjust the number of components as new data arrives. The variational inference method commonly used for DP-GMM in practice has different convergence properties.

**a. **Variational Inference (VI) in DP-GMM**:** - Variational inference (VI) is an optimization-based approach that approximates the posterior distribution by minimizing the **Kullback-Leibler (KL) divergence** between the true posterior and an approximate distribution. - **Convergence Guarantee**: The variational inference algorithm is guaranteed to converge to a **local optimum** of the variational objective (evidence lower bound, or ELBO). Like EM, it increases the ELBO monotonically at each iteration, but the solution is not guaranteed to be globally optimal.

- **Dirichlet Process Convergence**: The Dirichlet Process itself converges in the sense that the number of components will stabilize as more data is observed. However, convergence to the true underlying distribution depends on the amount of data and the correctness of the model assumptions.

**b. **MCMC-based DP-GMM**:** For **MCMC-based inference** in DP-GMM (e.g., Gibbs sampling): - **Convergence Guarantee**: MCMC methods for DP-GMM have well-known asymptotic guarantees. Given enough iterations, MCMC will sample from the true posterior distribution (including the true number of components), but in practice, convergence can be slow, and there are no guarantees on the rate of convergence.

**Summary of Convergence Guarantees:**

**Practical Considerations:** - **Incremental GMM** is useful for nonstationary or streaming data but lacks formal convergence guarantees. - **Online EM** is widely used in online or streaming contexts with convergence guarantees under proper conditions, but only to a local optimum. - **DP-GMM** with variational inference

provides a flexible and principled approach, but it also converges to local optima of the variational objective. MCMC methods for DP-GMM offer better posterior approximations but come with slower convergence.

In summary, **Online EM** provides the strongest mathematical guarantees on convergence (to a local optimum), while DP-GMM offers flexibility with some guarantees via variational inference or MCMC methods.