Unsupervised Phoneme Discovery Using Markov Models

Programming and Synthesis 2

March 2025

This document introduces the process of unsupervised phoneme discovery using Markov Models and clustering techniques. It is designed for Master's students in Electronic Music Composition as part of the "Programming and Synthesis 2" course. The focus is on exploring creative applications for audio synthesis, machine learning, and experimental sound design.

1 Introduction

In experimental music and sound design, machine learning can provide tools for discovering patterns in sound. One such application is the **unsupervised discovery of phonemes**, the fundamental sound units of speech. This text introduces a method using **Markov Models and clustering** to segment and analyze phoneme-like structures in an unknown language, useful for both linguistic exploration and generative sound synthesis.

1.1 Creative Use Cases

- Algorithmic Composition: Extract phoneme-like structures and use them the synthesis of new sounds.
- **Timbre Morphing:** Apply discovered phoneme transitions to resynthesize vocal-like textures. Or if not dealing with speech audio, to resynthesize audio that mimics the structures of the input.
- Live Performance Tools: Build an interactive system where Markovgenerated sounds influence synthesis in real time.
- **Sound Sculpting:** Use phoneme sequences as control signals for granularlike synthesis.
- Other?: Discuss your ideas for the creative use or misuse of this technique.

2 Conceptual Overview

The process of **unsupervised phoneme discovery** is an attempt to extract the fundamental building blocks of speech without any prior linguistic knowledge. Instead of predefined phoneme categories, we use machine learning techniques to identify patterns in the sound and represent them as probabilistic states in a Markov Model.

How to use!!: you will use this python script, create its necessary virtual environment to install required libraries and make sure you have a *train_speech.wav* and *test_speech.wav* audio files that you recorded previously, in the same folder where the script is. These audio files need to be at 16 kHz sampling rate.

The steps of our code are as follows:

- Audio Framing: Input is a continuous waveform. To analyze it meaningfully, we divide it into short, overlapping frames. Each frame is a snapshot of the sound over a short period (e.g., 25ms).
- **Spectral Feature Extraction**: The Fast Fourier Transform (FFT) is applied to each frame to extract spectral information, giving us a frequency-domain representation of each frame.
- Clustering with K-means: We use K-means clustering to group frames with similar spectral features into discrete categories, effectively forming a "codebook" of fundamental speech components. In our case we are arbitrarily choosing the number of clusters k to be 50.
- Markov State Assignment: Each frame is assigned to the nearest cluster center, forming a sequence of states that represent the speech.
- Markov Transition Modeling: By analyzing how states(clusters) transition from one to another, we compute a Markov Transition Matrix that captures the statistical relationships between different phoneme-like units.
- Markov-Based Synthesis: Using the learned transition probabilities, we generate new speech-like sounds by probabilistically selecting frames based on their transitions.

This approach is powerful because it does not require labeled phoneme data. Instead, it lets the structure of speech emerge naturally, making it useful for both scientific analysis and artistic exploration.

3 Implementation

3.1 Main Function: Running the Full Pipeline

The following function integrates all the steps into a complete workflow:

```
def main_pipeline(train_audio, test_audio, k_clusters=50):
       ""Full pipeline: From audio to Markov matrix and synthesis
2
      using k-means clustering."""
3
      # 1. Frame the training audio and extract FFT features
4
      train_frames, sr = frame_audio(train_audio)
5
6
      train_features = compute_fft(train_frames)
7
      # 2. Create the codebook using K-means.
8
      kmeans_model = build_codebook(train_features, k=k_clusters)
9
10
      # 3. Process the test audio and assign states
      test_frames, _ = frame_audio(test_audio)
12
      test_features = compute_fft(test_frames)
13
14
      assigned_states = assign_to_codebook(test_features,
      kmeans_model)
16
17
      # 4. Compute and visualize the Markov transition matrix
      markov_matrix = compute_markov_matrix(assigned_states,
18
      k clusters)
      visualize_markov_matrix(markov_matrix)
19
20
      # 5. Display sample state sequence and matrix
21
      print(f"\nSample State Sequence: {assigned_states[:20]}")
      print(f"\nMarkov Transition Matrix (shape {markov_matrix.shape
23
      }):\n", markov_matrix)
      # 6. Synthesize audio based on the Markov model
      synthesized_audio = synthesize_audio(test_frames,
26
      assigned_states, markov_matrix, length=2000)
      save_synthesized_audio(synthesized_audio, sr)
27
```

Listing 1: Main pipeline for phoneme discovery and synthesis.

Explanation: The function executes the pipeline in a structured manner:

- 1. It extracts **FFT features** from a reference (training) audio file.
- 2. It uses **K-means clustering** to create a set of representative sound units with similar spectral features (categories). This is the **codebook**
- 3. It processes a separate (test) audio file, **assigning states** to each of its **framed snippets** based on the trained **codebook**.
- 4. It computes the **Markov Transition Matrix**, which models how these states transition over time.
- 5. It generates a **new audio sequence** by stochastically (probabilistically) selecting frames based on the learned transitions.

3.2 Step 1: Frame the Audio

```
import librosa
1
2
 def frame_audio(audio_file, frame_size_ms=25, hop_size_ms=10, sr
3
     =16000):
      """Frame the audio file into overlapping frames."""
4
     y, sr = librosa.load(audio_file, sr=sr)
5
     frame_size = int(sr * frame_size_ms / 1000)
6
     hop_length = int(sr * hop_size_ms / 1000)
7
     frames = librosa.util.frame(y, frame_length=frame_size,
8
     hop_length=hop_length).T
     return frames, sr
```

Listing 2: Framing audio into overlapping segments.

3.3 Step 2: Compute FFT Features

```
import numpy as np
def compute_fft(frames):
    """Compute the FFT magnitude spectrum for each frame."""
    fft_frames = np.abs(np.fft.rfft(frames, axis=1))
    return fft_frames
```

Listing 3: Computing FFT magnitude spectrum for each frame.

3.4 Step 3: Cluster FFT Features Using K-Means

```
1 from sklearn.cluster import KMeans
2
3 def build_codebook(features, k=50):
4 """Cluster the FFT features using K-means to create the
        codebook."""
5 kmeans = KMeans(n_clusters=k, random_state=0, max_iter=300,
        n_init=10)
6 kmeans.fit(features)
7 return kmeans
```

Listing 4: Clustering FFT features using K-means.

3.5 Step 4: Compute Markov Transition Matrix

```
1 def compute_markov_matrix(states, num_states):
2 """Compute the Markov state transition matrix."""
3 matrix = np.zeros((num_states, num_states))
4 for (i, j) in zip(states[:-1], states[1:]):
5 matrix[i, j] += 1
6 matrix /= matrix.sum(axis=1, keepdims=True) + 1e-10 #
Normalize
7 return matrix
```

Listing 5: Computing Markov state transition probabilities.

3.6 Step 5: Generate Audio Using Markov Model

```
def synthesize_audio(frames, assigned_states, markov_matrix, length
1
      =1000):
      """Synthesize audio by selecting frames based on Markov
2
      transitions."""
      synthesized_audio = []
3
      current_state = np.random.choice(np.arange(len(markov_matrix)))
4
      for _ in range(length):
5
          frame_indices = np.where(assigned_states == current_state)
6
      [0]
7
          if len(frame_indices) == 0:
              current_state = np.random.choice(np.arange(len(
8
      markov_matrix)))
9
              continue
          chosen_frame = frames[np.random.choice(frame_indices)]
10
          synthesized_audio.extend(chosen_frame)
11
          next_state_probs = markov_matrix[current_state]
12
          current_state = np.random.choice(np.arange(len(
      markov_matrix)), p=next_state_probs)
      return np.array(synthesized_audio)
14
```

Listing 6: Synthesizing new audio from Markov model.

3.7 Step 6: Visualize Markov Transition Matrix

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4
  def visualize_markov_matrix(markov_matrix):
      plt.figure(figsize=(10, 8))
5
      sns.heatmap(markov_matrix, cmap="viridis", square=True)
6
      plt.title("Markov State Transition Matrix")
7
      plt.xlabel("Next State")
8
      plt.ylabel("Current State")
9
      plt.show()
10
```

Listing 7: Visualizing the Markov state transition matrix.

3.8 Step 8: Save Synthesized Audio

```
import soundfile as sf
def save_synthesized_audio(audio_array, sr, output_path="
    synthesized_audio.wav"):
    sf.write(output_path, audio_array, sr)
    print(f"Synthesized audio saved to {output_path}")
```

Listing 8: Saving the synthesized audio to a file.

4 Future Challenges and Next Steps

This tutorial covers the fundamental steps for unsupervised phoneme discovery and Markov-based synthesis. However, there are several directions for further research and experimentation:

- Try other sources and explore: Instead of speech, try to use other sources of musical audio or any sound you can think of to understand what the script is doing.
- Use different parameters: Play with the parameters of the steps (size of snippets of audio, sampling rate, number of clusters, etc.)
- **Try Feedback:** One idea is to use the output audio as input for another run of the script. It should converge into digital noise.
- Use other features: Instead of computing FFTs, try with other methods, such as MEL Frequency Cepstral Coefficients (MFCC).
- Improving Phoneme Separation: Exploring different clustering methods such as Gaussian Mixture Models (GMMs) to achieve better phoneme separation.
- Explore different parameters and windowing in re-synthesis: In this example, we are concatenating snippets of audio without overlap. You could try to implement overlapping and windowing (i.e. Triangular, Hann, and other windows).
- Explore options for [re]synthesis: Right now we are simply concatenating snippets of audio. But audio could be re-synthesized using the spectrogram data directly. You could research how this is done.
- **Applying Deep Learning:** Using neural networks to refine phoneme representation and improve synthesis quality.
- Enhancing Temporal Structure: Investigating higher-order Markov models to capture longer phoneme transitions.
- Implement HMM when phonemes = group of clusters: As of now, we are identifying a phoneme with one of the clusters. But this means that a phoneme is assumed to be 25ms long. In relaity a phoneme is longer and made of several clusters put together. This is a different model that needs to find that relationship. It would use a training algorithm called Baum-Welch. This is of course already implemented in some libraries, is just a matter of knowing how to use it.
- Live Performance Integration: Using real-time Markov-based phoneme synthesis as an interactive performance tool.

5 Conclusion

This method allows us to extract meaningful phoneme-like units from speech and explore them creatively in our context of electronic music/sound manipulation. Potential applications include voice-based synthesis, stochastic sound generation, and AI-assisted vocal textures, conceptual artwork, sound design, pranking your neighbours, etc.

Further Reading

- [1] Aalto University. Speech processing. https://wiki.aalto.fi/display/ ITSP/Introduction+to+Speech+Processing, 2024. Accessed: March 5, 2025.
- [2] BuiltIn. Markov chains: A gentle introduction. https://builtin.com/ machine-learning/markov-chain, 2024. Accessed: March 5, 2025.
- [3] Antoine Caillon and Philippe Esling. Rave: A variational autoencoder for fast and high-quality neural audio synthesis. arXiv preprint arXiv:2111.05011, 2021.
- [4] Mousumi Malakar and Ravindra B Keskar. Progress of machine learning based automatic phoneme recognition and its prospect. *International Jour*nal of Speech Technology, 2023.
- [5] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, et al. Wavenet: A generative model for raw audio. arXiv preprint arXiv:1609.03499, 2016.
- [6] Ashish Vaswani, Noam Shazeer, Niki Parmar, et al. Attention is all you need. Advances in neural information processing systems, 30, 2017.

Appendices

Complete Script

```
1 import librosa
2 import numpy as np
3 from sklearn.cluster import KMeans
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import sys
7 import soundfile as sf
8
9
10 # ----- STEP 1: FRAME AUDIO
      -----
11 def frame_audio(audio_file, frame_size_ms=25, hop_size_ms=10, sr
      =16000):
      """Frame the audio file into overlapping frames."""
12
      y, sr = librosa.load(audio_file, sr=sr)
13
      frame_size = int(sr * frame_size_ms / 1000)
14
      hop_length = int(sr * hop_size_ms / 1000)
15
      frames = librosa.util.frame(y, frame_length=frame_size,
16
      hop_length=hop_length).T
     return frames, sr
17
18
19
   ----- STEP 2: COMPUTE FFT
20
21 def compute_fft(frames):
      """Compute the FFT magnitude spectrum for each frame."""
22
     fft_frames = np.abs(np.fft.rfft(frames, axis=1))
23
      return fft_frames
24
25
26
27 # ----- STEP 3: K-MEANS CLUSTERING
  def build_codebook(features, k=50):
28
      """Cluster the FFT features using K-means to create the
29
      codebook."""
     print(f"Clustering {features.shape[0]} frames into {k} clusters
30
      ...")
      kmeans = KMeans(n_clusters=k, random_state=0, max_iter=300,
31
     n_init=10)
32
     kmeans.fit(features)
     print("Codebook creation complete.")
33
34
     return kmeans
35
36
37 # ----- STEP 4: ASSIGN STATES
38 def assign_to_codebook(features, kmeans_model):
      """Assign each frame to the nearest codebook cluster center.""
39
      assigned_states = kmeans_model.predict(features)
40
41
      return assigned_states
42
43
44 # ----- STEP 5: MARKOV TRANSITION MATRIX
```

```
_____
  def compute_markov_matrix(states, num_states):
45
       """Compute the Markov state transition matrix."""
46
      matrix = np.zeros((num_states, num_states))
47
      for (i, j) in zip(states[:-1], states[1:]):
48
         matrix[i, j] += 1
49
      matrix /= matrix.sum(axis=1, keepdims=True) + 1e-10 #
50
      Normalize
      return matrix
52
53
          ----- STEP 6: VISUALIZE MARKOV MATRIX
54
  # ----
55 def visualize_markov_matrix(markov_matrix):
      """Visualize the Markov state transition matrix.""
56
      plt.figure(figsize=(10, 8))
57
58
      sns.heatmap(markov_matrix, cmap="viridis", square=True)
      plt.title("Markov State Transition Matrix")
59
      plt.xlabel("Next State")
60
      plt.ylabel("Current State")
61
      plt.show()
62
63
64
    ----- STEP 7: SYNTHESIZE AUDIO
65 #
  def synthesize_audio(frames, assigned_states, markov_matrix, length
66
      =1000):
      """Synthesize audio by probabilistically selecting frames based
67
      on the Markov matrix.'
      synthesized_audio = []
68
      current_state = np.random.choice(np.arange(len(markov_matrix)))
69
70
      print(len(markov_matrix))
      print(np.where(assigned_states == current_state)[0])
71
72
      for _ in range(length):
73
          frame_indices = np.where(assigned_states == current_state)
      [0]
          if len(frame_indices) == 0:
74
75
              current_state = np.random.choice(np.arange(len(
      markov matrix)))
              continue
76
77
78
          chosen_frame = frames[np.random.choice(frame_indices)]
79
          synthesized_audio.extend(chosen_frame)
80
          # Select next state based on transition probabilities
81
          next_state_probs = markov_matrix[current_state]
82
          current_state = np.random.choice(np.arange(len(
83
      markov_matrix)), p=next_state_probs)
84
      return np.array(synthesized_audio)
85
86
87
      ----- STEP 8: SAVE SYNTHESIZED AUDIO
88
89 def save_synthesized_audio(audio_array, sr, output_path="
      synthesized_audio.wav"):
      """Save the synthesized audio to a file."""
90
```

```
sf.write(output_path, audio_array, sr)
91
       print(f"Synthesized audio saved to {output_path}")
92
93
94
                        ---- MAIN PIPELINE ---
95
  #
   def main_pipeline(train_audio, test_audio, k_clusters=50):
96
       """Full pipeline: From audio to Markov matrix and synthesis
97
       using k-means clustering."""
98
99
       # 1. Frame the training audio and extract FFT features
       train_frames, sr = frame_audio(train_audio)
100
       train_features = compute_fft(train_frames)
102
       # 2. Create the codebook using K-means.
       kmeans_model = build_codebook(train_features, k=k_clusters)
104
105
106
       # ===========
107
       # 3. Process the test audio and assign states
108
       test_frames, _ = frame_audio(test_audio)
test_features = compute_fft(test_frames)
109
       assigned_states = assign_to_codebook(test_features,
       kmeans_model)
       # 4. Compute and visualize the Markov transition matrix
114
       markov_matrix = compute_markov_matrix(assigned_states,
115
       k_clusters)
       visualize_markov_matrix(markov_matrix)
116
       # 5. Display sample state sequence and matrix
118
119
       print(f"\nSample State Sequence: {assigned_states[:20]}")
       print(f"\nMarkov Transition Matrix (shape {markov_matrix.shape
120
       }):\n", markov_matrix)
121
       # 6. Synthesize audio based on the Markov model
       synthesized_audio = synthesize_audio(test_frames,
123
       assigned_states, markov_matrix, length=2000)
       save_synthesized_audio(synthesized_audio, sr)
124
125
126
     ----- RUN THE SCRIPT ------
127
  if __name__ == "__main__":
128
       # Provide two audio files (can be the same if testing)
129
       train_audio_file = "train_speech.wav" # Reference audio for
130
       building the codebook
       test_audio_file = "test_speech.wav"
                                                # Audio to analyze and
       synthesize from
132
       # Run the pipeline
133
       main_pipeline(train_audio_file, test_audio_file, k_clusters=50)
134
```

Listing 9: Full implementation of phoneme discovery and synthesis.