

# MUSIC ANALYSIS WITH A BAYESIAN MODEL

---

Lilia Ben Baccar  
MS-DS, ENSAE Paris  
lilia.benbaccar@ensae.fr

Erwan Rahis  
MS-DS, ENSAE Paris  
erwan.rahis@ensae.fr

January 2021

## Abstract

The analysis of music has always been an interesting topic for music theorists for helping in music teaching, analysis of human perception of sounds and for design of music search and organization tools. Moreover, a typical goal is to segment a given piece with the goal of interring interrelationships among motive and themes within the music. And analyzing sequential data has been a problem in statistical modeling for several years. The goal of this project is to use Bayesian inference to analyze a music piece. We consider a music piece as a sequential data and we want to segment the piece using a non-parametric Bayesian model. The music, to be represented as a discrete sequence of observations is processed and then modeled using a Hidden Markov Chain model. The segmentation is inferred by the model and compared to our music-theoretic analysis. Our method is mainly a combination of the papers by Lu Ren and al. [2] and Emily B. Fox and al. [1].

## 1 Bayesian model

To solve this problem, the approach taken is a non-parametric Bayesian model based on the *Hidden Markov Model* and a prior generated by a *Hierarchical Dirichlet Process* (HDP-HMM). Since the basic HDP-HMM tends to over-segment the audio data which creates redundant states and rapidly switched among them, we describe an augmented HDP-HMM that provides effective control over the switching rate: the sticky HDP-HMM, a bayesian nonparametric hidden Markov models with persistent states. So we decided to 'augment' the method of the paper [3] by applying the method of [1]. For this purpose, we want to segment a waveform into a set of time intervals with no *a priori* on the number of intervals. The sequence is modelled using the Hidden Markov Model with the states to categorise the interval. In other words, we want the piece to be divided into intervals and each interval will be labelled to a state that we will later on analyse. For example a music piece can be divided into parts that are : intro, chorus, pre-chorus, outro. Those parts can be easily identified by the human ear and brain (music-theoretic analysis). We will compare our analysis to the results of the model.

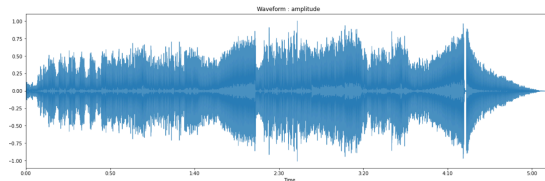


Figure 1: Audio waveform of the piece considered

### 1.1 Hidden Markov Model

The Hidden Markov Model used is described as follows : A discrete sequence of observed data  $y = \{y_t\}_{t=1}^T$  with  $y_t \in \{1, \dots, M\}$  and a corresponding hidden state sequence  $s = \{s_t\}_{t=1}^T$  with  $s_t \in \{1, \dots, I\}$ . The states will be analysed as in Section 1 (intro, chorus, etc.). Each observation is associated to a state. The discrete HMM model is represented by its parameters  $\theta = \{\pi, B, A\}$  such that :

- $\pi$  is a set of state transition probabilities, the probability to transition from one state to another. For two states  $\rho$  and  $\xi$ , the transition probability from  $\rho$  to  $\xi$  is  $\pi_{\rho\xi} = \mathbb{P}(s_{t+1} = \xi | s_t = \rho)$
- $B$  is a set of emission probabilities, the probability that an observation associated to a state  $\rho$  is equal to  $m$  :  $b_{\rho m} = \mathbb{P}(x_t = m | s_t = \rho)$

- $A$  is the set of initial state distribution, the probability that the first observation in time  $t = 1$  is equal to the state  $\rho$  :  $a_\rho = \mathbb{P}(s_1 = \rho)$

Because we will use a Hierarchical Dirichlet Process to model the parameters  $\theta$ , we will divide the sequence into  $J$  subsequences. This means that we will model  $J$  Hidden Markov Models. By taking the notation at the beginning of this section, we will have  $j \in \{1, \dots, J\}$  and  $y = \{y_j\}_{j=1}^J$  with  $y_j = \{y_{ji}\}_{i=1}^{N_j}$ . Note that every subsequence  $x_j$  has a length  $N_j$ . We allow the parameters to be different in each subsequence such that :

$$x_j \sim HMM(\theta_j)$$

The joint distribution given  $\theta$  is :

$$p(y|\theta) = \prod_{j=1}^J \left\{ \sum_{s_j} a_{s_j,1} \prod_{i=1}^{N_j-1} \pi_{s_j,i+1} \prod_{i=1}^{N_j} b_{s_j,i} y_{j,i} \right\}$$

## 1.2 Sticky Hierarchical Dirichlet Process

A Hierarchical Dirichlet Process is defined as a collection of Dirichlet Processes. Our collection of DP processes is the  $J$  parameters  $\theta_j$  linked to the HMM. This kind of setting allows for analysing how inter-related are the subsequences one to the others. The first process of the collection is denoted as :

$$G_0 \sim DP(\gamma, H)$$

with  $\gamma$  a concentration parameter for the Dirichlet distribution which is symmetrical and  $H$  a base probability measure on  $\Theta$ . The formal definition of a Dirichlet process is  $\forall (B_k)_{k=1,\dots,K}$  finite partition of  $\Theta$  :

$$(G_0(B_k))_{k=1,\dots,K} | \gamma, H \sim Dir((\gamma H(B_k))_{k=1,\dots,K})$$

We approach the Dirichlet Process with a stick-breaking process and we can write  $G_0$  as follows :

$$G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k}$$

with  $\delta_{\theta_k}$  an indicator function in  $\theta_k$  and  $\beta_k$  the mixture weights computed such that :

$$\begin{aligned} \beta_k &\sim Stick(\gamma) \\ \beta_k &= v_k \prod_{l=1}^{k-1} (1 - v_l) \\ v_k &\sim Beta(1, \gamma) \end{aligned}$$

with  $v_k$  the independent random variables following a Beta distribution. The process to obtain the probabilities  $\beta_k$  is pictured as breaking a unit-length stick because we start off with a stick of length 1 and at each step we break a portion of the what is remaining stick according to  $v_k$  and assign to  $\beta_k$  the piece broken. The smaller the  $\alpha$  the more concentrated is the distribution. The first process  $G_0$  (first subsequence) is then used as the base measure for the other processes (other subsequences) :

$$\begin{aligned} G_j &\sim DP(\alpha, G_0) \\ G_j &= \sum_{t=1}^{\infty} \tilde{\pi}_{jt} \delta_{\theta_{jt}^*} \\ \tilde{\pi}_j | \alpha &\sim Stick(\alpha) \\ \theta_{jt}^* | G_0 &\sim G_0 \\ \theta'_{ji} | G_j &\sim G_j \\ y_{ji} | \theta'_{ji} &\sim F(\theta'_{ji}) \end{aligned}$$

with  $j \in \{1, \dots, J\}$ ,  $i \in \{1, \dots, N_j\}$ . We know that multiple  $\theta_{jt}^*$  can take the same value  $\theta_k$ . We can write  $G_j$  as a function of unique states :

$$\begin{aligned} G_j &\sim \sum_{k=1}^{\infty} \pi_{jk} \delta_{\theta_k} \\ \pi_j | \alpha, \beta &\sim DP(\alpha, \beta) \\ \theta_k | H &\sim H \end{aligned}$$

with  $\pi_{jk}$  the distribution of the subsequence  $j$  over the state  $k$  with :

$$\pi_{jk} = \sum_{t|\theta_{jt}^*=\theta_k} \tilde{\pi}_{jt}$$

The model can be re written with an indicator variable  $z_{ji}$  for the state assigned to the observation  $y_{ji}$  :

$$\begin{aligned} G_j &\sim DP(\alpha, G_O) \\ z_{ji}|\pi_j &\sim \pi_j \\ y_{ji}|\{\theta_k\}, z_{ij} &\sim F(\theta_{z_{ji}}) \end{aligned}$$

With  $F(\cdot)$  the fixed distribution of the data sample. According to [1], we compute a sticky version of the model by adding a  $\kappa$  parameter that will slow down the transitions between states and avoid having unrealistically fast dynamics as results. The new transition distributions are modelled as follows :

$$\begin{aligned} \beta|\gamma &\sim Stick(\gamma) \\ \pi_j|\alpha, \kappa, \beta &\sim DP\left(\alpha + \kappa, \frac{\alpha\beta + \kappa\delta_j}{\alpha + \kappa}\right) \end{aligned}$$

By adding this parameter, we are increasing the probability of self-transition by an amount proportional to  $\kappa$ . If the probability that an observation has the same state as the one after is higher, the changing dynamics will slow down. That is why the  $\kappa$  parameter is added on the  $j^{th}$  component of  $\alpha\beta$ .

## 2 Computational method

We use the *Gibbs Block Sampler* [1] for the inference of the parameters of the sticky HDP-HMM. The first step is to sample  $\beta$  and  $\pi$ . The sampler is defined with a finite Dirichlet prior on  $\theta_j$  such that :

$$\begin{aligned} \beta|\gamma &\sim Dir(\gamma/L, \dots, \gamma/L) \\ \pi_j|\alpha, \beta &\sim Dir(\alpha\beta_1, \dots, \alpha\beta_L) \end{aligned}$$

This finite approximation is made because when  $L \rightarrow \infty$  then it converges to the HDP mixture model. The derived posterior distribution is :

$$\begin{aligned} \beta|\bar{n}, \gamma &\sim Dir(\gamma/L + \bar{m}_{\cdot 1}, \dots, \gamma/L + \bar{m}_{\cdot L}) \\ \pi_j|z_{1:T}, \alpha, \beta &\sim Dir(\alpha\beta + n_{j1}, \dots, \alpha\beta_j + \kappa + n_{jj}, \dots, \alpha\beta_L + n_{jL}) \end{aligned}$$

with  $n_{jk}$  the number of  $j$  to  $k$  transitions in the state sequence  $j$ . The choice of initialisation of the base measure is made regarding the paper [1] with :

$$\begin{aligned} y|z = k &\sim \mathcal{N}(\mu_k, \sigma_k^2) \\ \mu_k &\sim \mathcal{N}(0, s^2) \\ \sigma_k^2 &\sim InvGamma(a, b) \end{aligned} \tag{1}$$

After, choosing the hyper-parameters and base measure, the next step is to sample the hidden states  $z_t$ . This is done with a forward-backward algorithm to jointly sample  $z$  given  $y$  for the sticky HDP-HMM. The conditional distribution is :

$$p(z|y, \pi, \theta) = p(z_t|z_{t-1}, y, \pi, \theta)p(z_{t-1}|z_{t-2}, y, \pi, \theta) \dots p(z_2|z_1, y, \pi, \theta)p(z_1|y, \pi, \theta)$$

The algorithm starts by sampling  $z_1$  from  $p(z_1|y, \pi, \theta)$  then sample  $z_2$  from  $p(z_2|z_1, y, \pi, \theta)$  and so on. The conditional distribution is :

$$\begin{aligned} p(z_1|y, \pi, \theta) &\propto p(z_1)f(y_1|\theta_{z_1}) \sum_{z_2:T} \prod_t p(z_t|\pi_{z_{t-1}})f(y_t|\theta_{z_t}) \\ &\propto p(z_1)f(y_1|\theta_{z_1}) \sum_{z_2} p(z_2|\pi_{z_1})f(y_2|\theta_{z_2})m_{3,2}(z_2) \\ &\propto p(z_1)f(y_1|\theta_{z_1})m_{2,1}(z_1) \end{aligned}$$

where  $m_{t,t-1}(z_{t-1})$  is called the *backward message* passed from  $z_t$  to  $z_{t-1}$  and defined by :

$$m_{t,t-1}(z_{t-1}) \propto p(y_{t:T}|z_{t-1}, \pi, \theta)$$

In the algorithm the messages are first initialized to :

$$m_{T+1,T}(k) = 1$$

and then with a backward loop for each  $k \in \{1, \dots, L\}$  we compute the messages :

$$\begin{aligned} m_{t,t-1}(k) &= \sum_{j=1}^L \pi_k(j) \mathcal{N}(y_t, \mu_j, \Sigma_j) m_{t+1,t}(j) \\ &= \sum_{j=1}^L \pi_k(j) \frac{1}{\Sigma_j} \exp \left[ -\frac{1}{2} \left( \frac{y_t - \mu_j}{\Sigma_j} \right)^2 \right] m_{t+1,t}(j) \end{aligned}$$

Because of (1) we have  $\mathcal{N}(y_t, \mu_j, \Sigma_j) \propto \frac{1}{\Sigma_j} e^{-\frac{1}{2} \left( \frac{y_t - \mu_j}{\Sigma_j} \right)^2}$ . Then we derive the general distribution for  $z_t$  which is :

$$p(z_t | z_{t-1}, y, \pi, \theta) \propto p(z_t | \pi_{z_{t-1}}) f(y_t | \theta_{z_t}) m_{t+1,t}(z_t)$$

The algorithm computes the  $z_t$  with a forward loop with for each  $(j, k) \in \{1, \dots, L\}^2$  computes the probability of  $y_t | z_t = k$ :

$$f_k(y_t) = \pi_{z_{t-1}}(k) \mathcal{N}(y_t, \mu_k, \Sigma_k) m_{t+1,t}(j)$$

and then assigns a state  $z_t$  if the probability of having  $y_t$  conditional to the state  $k$  is inferior to the probability conditional to the state  $j$  :  $f_k(y_t) < f_j(y_t)$  then the state in  $t$  is updated  $z_t = j$  and we eventually update the transition matrix (N). The next step is to sample the auxiliary variables. They are introduced to simplify the inference algorithm. Those random variables are the following :

$$\begin{aligned} w_{jt} | \alpha, \kappa &\sim \text{Ber} \left( \frac{\kappa}{\alpha + \kappa} \right) \\ \bar{k}_{jt} | \beta &\sim \beta \\ k_{jt} | \bar{k}_{jt}, w_{jt} &= \begin{cases} \bar{k}_{jt}, & w_{jt} = 0 \\ j, & w_{jt} = 1 \end{cases} \end{aligned}$$

With  $k_{jt}$  the state actually assigned,  $\bar{k}_{jt}$  a considered state and  $w_{jt}$  an override variable that can override the considered state and be assigned instead. This will allow for the  $\kappa$  parameter to be taken into account for the "stickiness" of the model. We want at the end a higher probability of staying at the same state to avoid the fast state change dynamic.

We have to consider the number of observations that are transitioning from state  $j$  to  $k$  and for each transition  $n_{jk}$  the probability that an observation actually transitions is given by a Bernoulli random variable:

$$\text{Ber} \left( \frac{\alpha \beta_k + \kappa \delta(j, k)}{n + \alpha \beta_k + \kappa \delta(j, k)} \right)$$

For  $i = 1, \dots, n_{ij}$  we sample the random variable and increment the count  $m$  when the variable is equal to 1, then  $m_{jk}$  follows a Binomial distribution :

$$m_{jk} \sim \text{Binomial} \left( n_{jk}, \frac{\alpha \beta_k + \kappa \delta(j, k)}{n + \alpha \beta_k + \kappa \delta(j, k)} \right)$$

Then for each  $j \in \{1, \dots, L\}$  we sample the number of override variables  $w_j$ . in state  $j$ . As we draw  $m_{jj}$  sample of  $w_{jt}$  being a Bernoulli variable, then the sum  $w_j$ . is a binomial:

$$w_j. \sim \text{Binomial} \left( m_{jj}, \frac{\kappa}{\kappa + \alpha \beta} \right)$$

Then we compute  $\bar{m}_{jk}$  the number of observations that consider the state  $k$  and remove the override variable  $w_j$  from the  $m_{jj}$  counts. The next step is to compute the  $\beta$  from their distribution :

$$\beta^{(n)} \sim \text{Dir}(\gamma/L + \bar{m}_{.1}, \dots, \gamma/L + \bar{m}_{.L})$$

Then for each  $k \in \{1, \dots, L\}$  we sample a new transition distribution  $\pi_k$  following :

$$\begin{aligned} \pi_k &\sim \text{Dir}(\alpha \beta_1 + n_{k1}, \dots, \alpha \beta_k + \kappa + n_{kk}, \dots, \alpha \beta_L + n_{kL}) \\ \theta_k &\sim p(\theta | \lambda, \mathcal{Y}_k) \end{aligned}$$

By taking equations in (1) as base measures we have the parameters  $\theta_k = \{\mu_k, \Sigma_k\}$ . We placed a Gaussian prior over the mean of the distribution and an inverse-Gamma over the variance and at the end of the sampler we compute  $\mu$  and  $\sigma$  for the

### 3 Data processing

We consider an acoustic signal in a mp3 format loaded in Python using the librosa package. The sampling rate is 22.05kHz which means we have 22050 observations per second of audio signal. The acoustic signal is first loaded as a waveform or amplitude as shown in Figure 1. It represents the pressure change recorded by a recording instrument such as a microphone. We then divide the piece into 100ms contiguous frames. Each amplitude frame is processed by computing its Mel frequency cepstral coefficients (MFCCs). These coefficients are derived from the Melspectrogram that is a log-scale of pitches judged to be equally distanced to one another. The Mel spectrogram is the representation showing the mel frequencies according to time and intensity. The MFCCs are the most popular feature used in audio processing. We compute 40 columns of coefficients for each frame. To obtain a vector sequential data, we quantize each frame of coefficient. First, we *whiten* the columns so each one of them have a unit variance. Indeed, before running k-means algorithm, we need to rescale each feature dimension of the observations by the standard deviation. Then we use the k-means algorithm with 16 centroids to generate to adjust the classification of the observations into clusters and updates the clusters centroids until the position of the centroids is stable over successive iterations. This yields a code book mapping centroids to codes and vice-versa. After that, we pass them to the vector quantization algorithm which assigns codes from a code book to observations. The result is a 1D discrete sequential data representing the coefficients. After the piece is processed, the sequence is represented as a series of discrete observations. According to the paper, we should have the following figure which represents each observation by its cluster discretization.

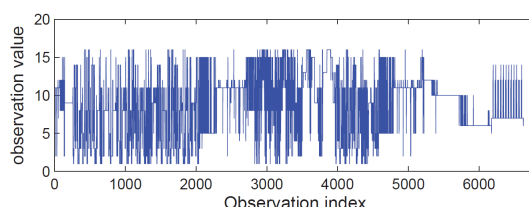


Figure 2: Sequence of code indices for the Beatles' music using a codebook of dimension  $M = 16$

Unfortunately, after a lot of attempts, we've never managed to find this plot. After the discretization using k-means and Vector Quantization, we find the following plot:

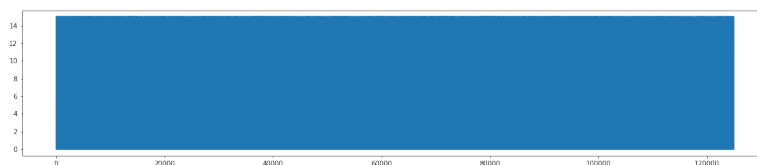


Figure 3: Our output

Indeed, the paper doesn't give enough details about the pre-processing part. We tried to use other papers from the same author, about the same topic, to see if he's giving some other details but we didn't succeed. So, we decided to use simulated data to apply our bayesian model. We found them online in the Github of a similar topic.

### 4 Results

Since our results are not interpretable, we propose to analyze the results of the research paper studied. As we said, after the piece is processed, the sequence is represented as a series of discrete observations as we showed in Figure 3. Then, we divided the piece into 88 subsequences of 75 observations. The music segmentation results using dHDP HMM model is shown in the following figure :

The first plot is the similarity matrix  $E(z'z)$ . It quantifies how inter-related any one subsequence of the music is to all others. The second one is the segmentation result on the Beatles audio waveform 1. We can see that the audio is decomposed into clear segments of various lengths and some of them are repeated. To analyse the second plot, we have to compare these segments with a music-theoric analysis of the same audio. This method will allowed us to compare the results of our model and the segmentation done by a person.

As we see, the segmentation is almost the same. Our music audio is divided into 10 parts with some redundant segments corresponding to the choruses. Each segment represents a specific part of the music: the verse of a singer, the instrumental part,...

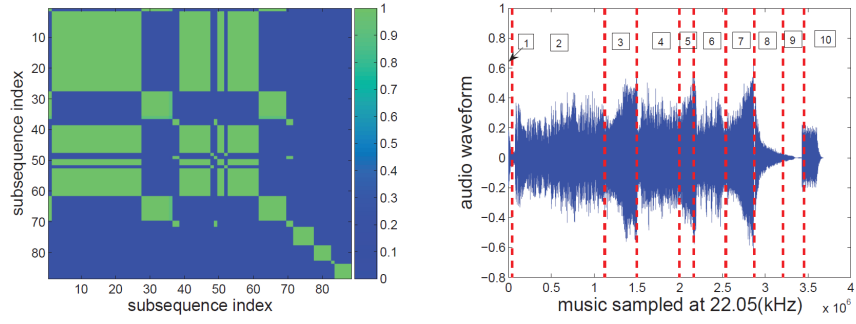


Figure 4: Paper results

segment index	subsequences included	music theory explanation
1	1 <sup>st</sup>	an instrumental accompaniment with some applause
2	2 <sup>nd</sup> ~ 27 <sup>th</sup>	three verses sung by Lennon
3	28 <sup>th</sup> ~ 36 <sup>th</sup>	an orchestral crescendo continues
4	37 <sup>th</sup> ~ 47 <sup>th</sup>	an interlude ('Woke up....') sung by McCartney
5	48 <sup>th</sup> ~ 52 <sup>nd</sup>	a short transition
6	53 <sup>rd</sup> ~ 61 <sup>st</sup>	a verse part sung by Lennon
7	62 <sup>nd</sup> ~ 69 <sup>th</sup>	the same orchestral crescendo as the third part
8	70 <sup>th</sup> ~ 77 <sup>th</sup>	the "famous" final chords played on three different pianos
9	78 <sup>th</sup> ~ 82 <sup>nd</sup>	an almost quiet part
10	83 <sup>rd</sup> ~ 88 <sup>th</sup>	the famous "studio chatter" part

Figure 5: Paper results

## References

- [1] Emily B. Fox, Erik B. Sudderth, Michael I. Jordan, and Alan S. Willsky. A sticky HDP-HMM with application to speaker diarization. *The Annals of Applied Statistics*, 5(2A):1020--1056, June 2011.
- [2] Lu Ren, David B. Dunson, Scott Lindroth, and Lawrence Carin. Music analysis with a Bayesian dynamic model. *IEEE*, 2009.
- [3] Lu Ren, David Dunson, Scott Lindroth, and Lawrence Carin. Dynamic Nonparametric Bayesian Models for Analysis of Music. *Journal of the American Statistical Association*, 105(490):458--472, June 2010.
- [4] Satyanand Singh and E.G. Rajan. Vector Quantization Approach for Speaker Recognition using MFCC and Inverted MFCC. *International Journal of Computer Applications*, 17(1):1--7, March 2011.
- [5] IEEE Signal Processing Society, editor. *Proceedings / 2009 IEEE International Conference on Acoustics, Speech, and Signal Processing: April 19 - 24, 2009, Taipei International Convention Center, Taipei, Taiwan*. IEEE, Piscataway, NJ, 2009. Meeting Name: IEEE International Conference on Acoustics, Speech, and Signal Processing.