
Online Speaker Diarization with interactive learning

Anirudh Garg
anirudhg@iitk.ac.in
170117

Anubhav Satpathy
asatpath@iitk.ac.in
170136

Nitish Vikas Deshpande
nitishvd@iitk.ac.in
17807450

Abstract

Speaker diarization is the task to partition an audio stream into homogeneous segments according to the speaker identity. A layman's way to put it would be "Who spoke when". It is observed that some state-of-the-art speaker diarization systems require really large datasets to train the clustering modules which might not be easily available everywhere. Here, the method of learning continually can be employed i.e., online learning. Online learning is a problem where data becomes available in a sequential order and later used to update the best predictor for future data or reward associated with the data features. The only way sometimes, in which the online learning agent can learn from the past experience is the feedback in terms of rewards approach. This online learning problem is particularly important in the field of sequential decision making. In sequential decision making, the best possible action to perform at each step to maximize the cumulative reward over time is chosen by the agent. It is important to obtain an optimal balance between the exploration of new actions and the exploitation of the possible rewards generated from known previous actions.

1 Baseline Method

We formulate the online speaker diarization as a contextual-bandit problem similar to the online semi-supervised learning method in [2]. In a bandit problem, each arm of the bandit corresponds to a certain distribution of probability for the rewards and in each round, a particular arm is chosen by the agent. Reward and updates are received by the agent based on this move. A more refined version of this is known as the contextual bandit [1], where the results are obtained given the context based on a relationship between the feature vectors and the context.

1.1 A Multi-Armed Bandit Formulation

Formally, a contextual-bandit algorithm proceeds in discrete trials (t). In a trial t :

- A current user u_t and a set \mathcal{A}_t of arms and actions is considered by the algorithm along with the corresponding feature vectors $x_{t,a}$, where $a \in \mathcal{A}_t$. Through the vector $x_{t,a}$, we obtain the information about both the user u_t and the chosen arm a . This vector will be referred to as the *context* in future.
- Analysing the payoffs in previous trials, an arm $a_t \in \mathcal{A}_t$ is chosen by A, which receives a payoff of r_{t,a_t} . The expectation of this payoff is dependant on both the user u_t and the arm a_t chosen by the user.
- The strategy to select the arms is improved by the new observation, $(x_{t,a}, u_t, r_{t,a})$. An important thing to note is that for the unchosen arms, $a \neq a_t$, we will receive no feedback i.e., no payoff $r_{t,a}$.

In the event of total N trials, the payoff is described by $\sum_{t=1}^N r_{t,a_t}$, thus the expected payoff will be given by $\mathbb{E}\left[\sum_{t=1}^N r_{t,a_t^*}\right]$, where a_t^* denotes the arm which has the maximum payoff at the trial 't'. Our goal is to maximize the expected payoff that is described above.

1.2 The K-armed Bandit

K-armed bandit is a specific case of the general contextual bandit problem in which the user u_t is the same for all the trials and also the corresponding set of arms \mathcal{A}_t remains unchanged and contains K arms for all the trials. As the arm set and context are invaried in each trial, they make no difference whatsoever to the bandit algorithm, thus it can be called as the context free bandit algorithm.

1.3 The LinUCB Algorithm

If we are given a parametric form of payoff function, an efficient method to compute the confidence interval, from the data, of the parameters, with which we can compute a UCB (Upper Confidence Bound) of the estimated arm payoff, is the LinUCB algorithm. The LinUCB algorithm results in a closed form solution when the payoff model is linear. Assuming that the expected payoff of an arm a is linear in its d -dimensional feature $x_{t,a}$ and employing an unknown vector θ_a^* acting as the coefficient, we get,

$$\mathbb{E}[r_{t,a_t} | \mathbf{x}_{t,a}] = \mathbf{x}_{t,a}^T \theta_a^* \quad (1)$$

At the t^{th} trial, let \mathbf{D}_a be a matrix of dimension $m \times d$. In this matrix, the rows correspond to the m contexts which were observed previously for the arm a . $c_a \in \mathbb{R}^m$ is the corresponding response vector. On implementing ridge regression to the data (\mathbf{D}_a, c_a) we get an estimate of the coefficients:

$$\hat{\theta} = (\mathbf{D}_a^T \mathbf{D}_a + \mathbf{I}_d)^{-1} \mathbf{D}_a^T c_a \quad (2)$$

where \mathbf{I}_d is the $d \times d$ identity matrix. It can be shown that the payoff for arm a can have a reasonable UCB such that:

$$| \mathbf{x}_{t,a}^T \hat{\theta}_a - \mathbb{E}[r_{t,a_t} | \mathbf{x}_{t,a}] | \leq \alpha \sqrt{\mathbf{x}_{t,a}^T (\mathbf{D}_a^T \mathbf{D}_a + \mathbf{I}_d)^{-1} \mathbf{x}_{t,a}} \quad (3)$$

where $\mathbf{x}_{t,a} \in \mathbb{R}^d$ and $\alpha = 1 + \sqrt{\ln(2/\delta)/2}$. Thus, a_t can be reasonably estimated as:

$$a_t = \operatorname{argmax}_{a \in \mathcal{A}_t} (\mathbf{x}_{t,a}^T \theta_a + \alpha \sqrt{\mathbf{x}_{t,a}^T \mathbf{A}_a^{-1} \mathbf{x}_{t,a}}) \quad (4)$$

where $\mathbf{A}_a = \mathbf{D}_a^T \mathbf{D}_a + \mathbf{I}_d$. A comprehensive description of the LinUCB algorithm is given by Algorithm 1. It is useful to keep in mind that LinUCB always chooses the arm with the highest UCB. We see that the computational complexity of this algorithm is at most cubic in number of features and is linear in the number of arms. Computation complexity can be decreased further by updating \mathbf{A}_{a_t} in every step (which takes $O(d^2)$ time). The algorithm can be employed for a dynamic arm set too and gives efficient results as long as we make sure that the size of \mathcal{A}_t is not too large.

Algorithm 1 LinUCB

```
1: Initialize  $c_t \in \mathbb{R}_+$ ,  $\mathbf{A}_a \leftarrow \mathbf{I}_d$ ,  $\mathbf{b}_a \leftarrow \mathbf{0}_{d \times 1} \forall a \in \mathcal{A}_t$ 
2: for  $t = 1, 2, 3, \dots, T$  do
3:   Observe features  $\mathbf{x}_t \in \mathbb{R}^d$ 
4:   for all  $a \in \mathcal{A}_t$  do
5:      $\hat{\theta}_a \leftarrow \mathbf{A}_a^{-1} \mathbf{b}_a$ 
6:      $p_{t,a} \leftarrow \hat{\theta}_a^T \mathbf{x}_t + c_t \sqrt{\mathbf{x}_t^T \mathbf{A}_a^{-1} \mathbf{x}_t}$ 
7:   end for
8:   Choose arm  $a_t = \operatorname{argmax}_{a \in \mathcal{A}_t} p_{t,a}$ 
9:   Observe feedback  $r_{a_t,t}$ 
10:   $\mathbf{A}_{a_t} \leftarrow \mathbf{A}_{a_t} + \mathbf{x}_t \mathbf{x}_t^T$ 
11:   $\mathbf{b}_{a_t} \leftarrow \mathbf{b}_{a_t} + r_{a_t,t} \mathbf{x}_t$ 
12: end for
```

2 Dataset and Implementation

We used VoxCeleb [3], which is a large scale speaker recognition dataset to generate 3 different kinds of data which correspond to 5, 10, 15 speakers to mimic the real world conversations. Three different kinds of reward streams were used with epiReward being 0.01, 0.1, 0.5. Thus a total of 9 learning environments were analysed. To evaluate the performance, we measured the accuracy, which is the ratio of the correct identification of speakers and the total number of time steps or more accurately the ratio of total reward and the total number of time steps. We used Mel-frequency cepstral coefficients (MFCC) as the feature embedding. Using the spectrogram of the data, the MFCC creates a feature vector for the data that is further used in our computations.

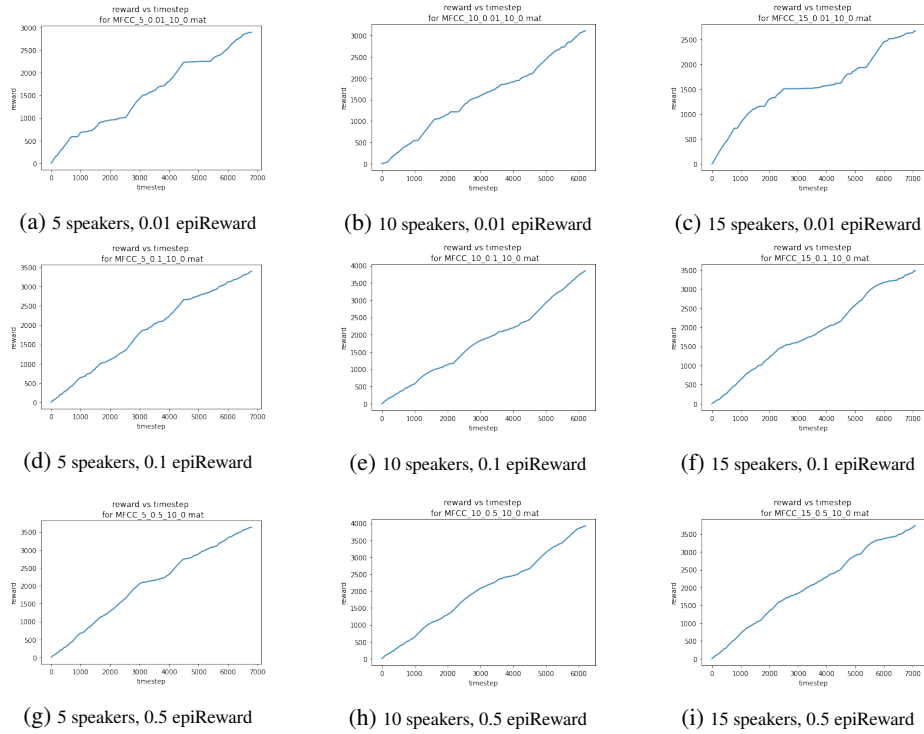


Figure 1: Plots for reward vs timestep for different number of speakers and epiRewards

3 Results

In figure 1, accuracy is calculated for different number of speakers (5, 10, 15) and different values of epiRewards (0.01, 0.1, 0.5). Sparsified reward streams with a revealing probability are known as the epiRewards. We can clearly observe from the table below that for higher epiReward, model achieves better accuracy. The model performs best for the case with 10 speakers.

Accuracy			
	5 speakers	10 speakers	15 speakers
LinUCB (epiReward = 0.01)	0.4253	0.5018	0.3770
LinUCB (epiReward = 0.1)	0.4994	0.6213	0.4896
LinUCB (epiReward = 0.5)	0.5346	0.6334	0.5258

References

- [1] Lihong Li et al. “A contextual-bandit approach to personalized news article recommendation”. In: *Proceedings of the 19th international conference on World wide web*. 2010, pp. 661–670.
- [2] Baihan Lin. “Online semi-supervised learning in contextual bandits with episodic reward”. In: *Australasian Joint Conference on Artificial Intelligence*. Springer. 2020, pp. 407–419.
- [3] Arsha Nagrani, Joon Son Chung, and Andrew Zisserman. “Voxceleb: a large-scale speaker identification dataset”. In: *arXiv preprint arXiv:1706.08612* (2017).