

# Exploring better ways to segment lecture videos based on topic transition

Project theme: Intelligent Learning Platform

Milad Parvaneh

*Email address:* `sparva2@illinois.edu`

October 23, 2022

## Team members

This is a one-member team. My name is Milad Parvaneh, and my NetID is sparva2. The captain is Milad Parvaneh.

## Topic selection

The topic I chose is: “Exploring better ways to segment lecture videos based on topic transition”. The idea is to improve an existing in-progress project, Smartmoocs, in which the goal is to improve students’ learning experience in MOOCs. The lecture videos are currently divided into 1-minute segments. Here, I aim to explore alternative ways using which we can segment the lecture videos according to the various topics presented in them. To achieve this, we would need to identify the transition points of topic changes and then make the segmentation based on these points.

The subject of the project is related to the second theme, Intelligent Learning Platform, as it tries to improve the learning experience within MOOCs. Furthermore, different methods for topic extraction can also be viewed as a subject of CS410, Text Information Systems course, since it can be subcategorized under the text mining and analytics title.

## Datasets, algorithms and techniques

The primary focus of this project will be on the lecture videos of the course CS410. For this course, lecture transcripts and time stamps are available on the Coursera platform. These transcripts will be used to extract various topics within a given lecture as well as their corresponding topic change points. Then, the time stamps associated with these transition points can be used to segment the videos.

Probabilistic topic models [Blei, 2012] will be among the methods considered for this project. This includes mixture models. Also, the Expectation-Maximization

(EM) algorithm will be employed, a practical algorithm for computing the maximum likelihood estimate of mixture models. Probabilistic Latent Semantic Analysis (PLSA) [Mei et al., 2007; Lu et al., 2011] is a model that can be used to mine multiple topics from text. It can be considered a mixture model with  $k$  unigram language models in which  $k$  represents the number of topics. Therefore, this method can be an ideal candidate for our video segmentation task.

Latent Dirichlet Allocation (LDA) is an extension of PLSA that can incorporate prior knowledge to extract topics. I noticed that there is a dataset for concept summary of lecture videos which provides a source for possible topics for each lecture. Using such a dataset as prior knowledge can potentially help us with improving the performance of the topic extraction task.

## Functioning expectation and programming language

PLSA models have been broadly used in practice as a primary topic model and are known to deliver acceptable results for most applications. Therefore, it seems it would be a suitable choice for this application. Also, since the lecture videos are currently divided into 1-minute segments, any contribution from a working topic model would be an improvement.

In this project, I will be using python as my main programming language.

## Workload

Table 1 lists the main tasks to be performed along with the estimated time cost for each task.

Table 1: Workload description.

Main tasks	Estimated time cost (Hrs)
Exploring mixture models including PLSA and LDA	3
Preparing input data (lecture video transcripts)	3
Implementing PLSA	5
Implementing LDA	5
Identify the transition points of topic changes	2
Analyzing and comparing results	2

# References

- D. M. Blei. Probabilistic topic models. *Communications of the ACM*, 55(4):77–84, 2012.
- Y. Lu, Q. Mei, and C. Zhai. Investigating task performance of probabilistic topic models: an empirical study of plsa and lda. *Information Retrieval*, 14(2):178–203, 2011.
- Q. Mei, X. Shen, and C. Zhai. Automatic labeling of multinomial topic models. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 490–499, 2007.