

ASSIGNMENT 6:

a) We used Gensim api in python for topic modeling.

b) The questions we thought of investigating are:

- How are the topics distributed for Election Candidates conversation over the debate?
- After generating the topic distribution how accurate are we?
- Do the topics make sense for the given Candidate conversations?

Data Source : <https://www.kaggle.com/mrisdal/2016-us-presidential-debates>

Data Description : Data is a .csv file that has the the conversation transcript for each Election candidate.

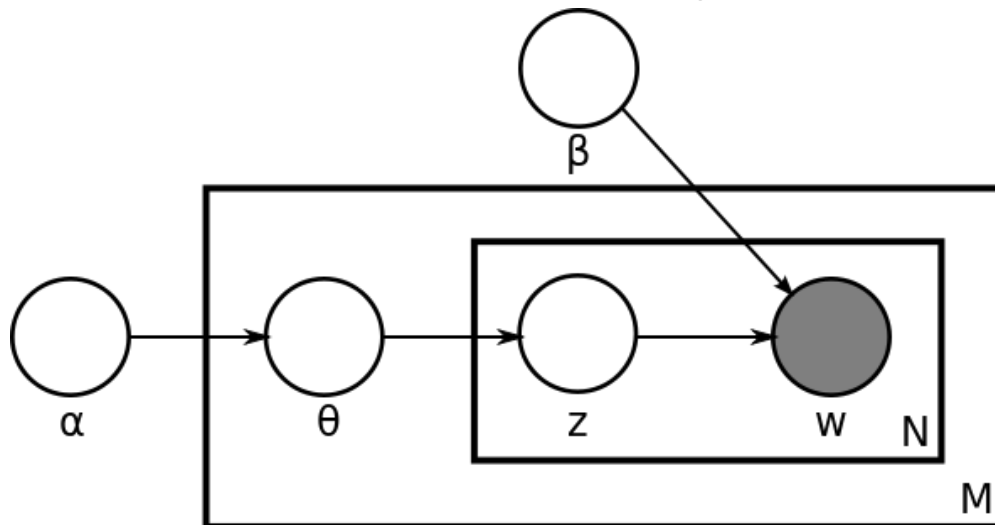
c)

We used LDA algorithm for Topic modeling.

In natural language processing, **Latent Dirichlet allocation (LDA)** is a generative statistical model that allows sets of observations to be explained by unobserved groups that explain why some parts of the data are similar. For example, if observations are words collected into documents, it posits that each document is a mixture of a small number of topics and that each word's creation is attributable to one of the document's topics.

Reference: https://en.wikipedia.org/wiki/Latent_Dirichlet_allocation,
<https://www.quora.com/What-is-a-good-explanation-of-Latent-Dirichlet-Allocation>

Plate Notation representing the LDA model



Thus, it represents documents as mixtures of topics that consist of certain words with certain probabilities.

It assumes that the documents are created in the following way:

- Decide on the number of words that the document will have (by Poisson Distribution)
- Choose a topic mixture for the document

- Generate each word in the document by: first picking the topic and then using the topic to generate the word itself.

LDA then tries to backtrack from the documents to find a set of topics that are likely to have generated the collection.

After generating the topics using gensim package in python for each election candidate, some topics generated did make sense, however topics with high words will **mr., go, think, just, want** had lower interpretation of the allocated topic. While word distribution for the topics such as **Country, Tax, Pay, American, People, Enforce, Terrorist** did make sense. One thing we denoted was, Clinton has higher probability for topic that has higher probability of word Trump, and V.i.Z. We generated 50 topics with 30 words that describes the topic. Hence, it does help us to understand the emotions, personality of each election candidate. If we were not computer savvy, using the api gensim functionality is difficult as it requires knowledge of object oriented programming with the way LDA is encompassed into the api. Also tuning the parameters for the training requires exponential amount of time with respect to increase training iterations and number of topics generated.

We have multinomial distribution over all the assigned documents for a respective document(candidate)

```
Election Candidate: Clinton : [(8, 0.078684132823948679), (38, 0.038178881074114327), (38, 0.00983888780189988), (44, 0.000203441888811844)]
Election Candidate: Raddatz : [(31, 0.31981188147114273), (38, 0.07889984718878817)]
Election Candidate: Wallace : [(8, 0.45506643912064904), (31, 0.37639547758851477), (38, 0.1058841378444732), (48, 0.043744749134823874)]
Election Candidate: Holt : [(8, 0.40816389489485987), (31, 0.410861347138848891), (38, 0.172188129711889448)]
Election Candidate: QUESTION : [(15, 0.99383647798741914)]
Election Candidate: Quijano : [(8, 0.53284795865556276), (31, 0.18844178888888888), (37, 0.31983198319831983), (48, 0.14188188188188188)]
Election Candidate: Trump : [(8, 0.93885653528160851), (31, 0.01855321555144447), (38, 0.844213413811384444)]
Election Candidate: Pence : [(37, 0.99928183468711796)]
Election Candidate: Audience : [(32, 0.97351351351351212)]
Election Candidate: CANDIDATES : [(38, 0.97958333333333333)]
Election Candidate: Cooper : [(31, 0.37651486398822442), (38, 0.41199178888888888)]
Election Candidate: Keane : [(37, 0.098488484229888988), (48, 0.08133888888888888)]
```

We also generate words per document with Topic ID as follows:

```
Topic: 0 word distribution is: (u'you.' : u' 0.013', u'trump.' : u' 0.008', u'question.' : u' 0.011', u'two' : u' 0.038', u'clinton.' : u' 0.034', u'audi' : u' 0.007', u'ask' : u' 0.014', u'you.' : u' 0.013', u'want' : u' 0.021', u'question' : u' 0.029', u'thank' : u' 0.022', u'trump,' : u' 0.037', u'clinton.' : u' 0.018', u'secretari' : u' 0.068', u'—' : u' 0.007', u'minutes.' : u' 0.010', u'please' : u' 0.010', u'debat' : u' 0.010', u'mr.' : u' 0.065', u'go' : u' 0.011', u'presidenti' : u' 0.008', u'next' : u' 0.012', u'let' : u' 0.016', u'respond.' : u' 0.007', u'american' : u' 0.018', u'final' : u' 0.007', u'move' : u' 0.016', u'clinton?' : u' 0.009', u'minut' : u' 0.013', u'issu' : u' 0.008')
```

Please run the program as follows placing debate.csv in the same folder as topic_model.py

```
python topic_model.py debate.csv
```