

Assignment-4 Report

Topic Modeling

Samreet Singh Dhami
2018CSB1120

IIT Ropar
Rupnagar

Abstract

The report is designed to develop a feel for the way topic modelling works, the connection to the human meanings of documents, and common ways of handling a time dimension. Analysis includes the State of the Union speeches corpus, and report on how the subjects have shifted over time in relation to historical events. Topic modeling has been done using two of the common methods, Latent Semantic Analysis(LSI) and Latent Dirichlet Allocation(LDA). Visualization on data modeling has been done using different techniques and observations are made on the results. At last, a new document set, collection of AP wire stories was introduced and results of two documents were compared.

1 Pre-processing and generating tf-idf score

Topic Modeling is a technique to extract the hidden topics from large volumes of text. LSI and LDA are two main algorithms in the gensim package of python to extract good quality topics. To achieve this, pre-processing of the content is necessary. Tokenization of the content and removing the stopwords is also very necessary. Once all these important steps are done, gensim is able to generate tf-idf scores.

2 Latent Semantic Indexing(LSI) on the data

Now it's time to apply Latent Semantic Indexing (LSI) to the tf-idf vectors. The major challenge arises is to determine number of topics. Once the number of topics are determined, we can import lsi model and apply it to print topics. Topics are printed in well-defined way. For example Topic 0 is represented as $0.072 \times \text{"upon"} + 0.062 \times \text{"tonight"} + 0.059 \times \text{"program"} + 0.058 \times \text{"mexico"} + 0.057 \times \text{"economic"} + 0.052 \times \text{"treaty"} + 0.051 \times \text{"subject"} + 0.050 \times \text{"help"} + 0.049 \times \text{"budget"} + 0.049 \times \text{"americans"}$

It means the top 10 keywords that contribute to this topic are: 'upon', 'tonight', 'program', .. and so on and the weight of 'tonight' on topic 0 is 0.062. Let us consider the first ten topics in the lsi model for further analysis.

1. The first topic has keywords program, Mexico, Economic, Treaty. This topic is dominant in initial years. If we deeply read the content, in the first few years after 1790 upto 1816 this topic does not hold good as in these years mostly issues related to Europe and tribes from India and other parts of the world are discussed. But in the years

```
[ (0, 0.006991717118536367), (1, 0.11238268667584739), (2, 0.01072822517259915), (3, 0.062293768326600024), (4, 0.05936337141487709), (5, 0.004925858332705161), (5, 0.023344416443247366), (6, 0.004058305176631579), (8, 0.020403093234589675), (13, 0.0055770015155953345), (0, 0.005311008274933625), (2, 0.01629862641781277), (6, 0.0014585405506684912), (8, 0.02199839989651333), (10, 0.03778386083862339), (0, 0.0059193203541003546), (6, 0.0032511976340816354), (7, 0.01854186028915431), (8, 0.012259025423878628), (15, 0.007513064659953558), (0, 0.003610550817851497), (5, 0.017110965138726713), (6, 0.0009915507143285767), (8, 0.007477518615811034), (9, 0.009619229904059107), (0, 0.0049857592799054905), (8, 0.010325602622159412), (9, 0.013283062286262718), (10, 0.023646637958361243), (11, 0.0275926386286267288), (0, 0.01335042411227495), (5, 0.02108891903806155), (6, 0.0012221241823501216), (7, 0.013939759064810268), (8, 0.009216327710019965), (3, 0.023907859178310897), (6, 0.002947687607378017), (9, 0.028596101410622396), (10, 0.01272676515597039), (11, 0.014850526847574368), (0, 0.004740846070956307), (3, 0.02111964498915876), (6, 0.003905877687068892), (8, 0.009818382692244325), (10, 0.011242527402801948), (6, 0.00358103153489415), (13, 0.00783168426040194), (15, 0.008275264840715354), (29, 0.0038138652135741953), (30, 0.002617408629425127), (13, 0.007819295788240576), (20, 0.015949309658042657), (27, 0.0035478677626071737), (42, 0.02299460952125571), (46, 0.0545158832078651), (0, 0.012724417127693979), (2, 0.0065081967620449705), (9, 0.011300132756426166), (10, 0.01005830366575092), (11, 0.023473538923473675), (0, 0.00517134262705343), (2, 0.031739917499661145), (6, 0.004260539009558743), (9, 0.013777444603237807), (10, 0.012263371107620007), (0, 0.005251693937586234), (6, 0.005769005184055994), (11, 0.014532199096546227), (15, 0.01333136703778494), (20, 0.024256185031136957), (6, 0.006132336309554068), (8, 0.011561349876126788), (9, 0.014872752335858483), (11, 0.015447434927542004), (13, 0.006320381428456844), (0, 0.008644795883484928), (2, 0.006632370506188272), (6, 0.0011870425252229898), (7, 0.04061883491120717), (8, 0.017903537261076297), (0, 0.004312958156733443), (6, 0.003553350643246007), (7, 0.013510043293195597), (8, 0.00893221940192548), (10, 0.01027827182278212), (5, 0.023379258097285915), (9, 0.013143061002160606), (11, 0.027301816771343673), (13, 0.005855325217232992), (20, 0.01139256448965873), (2, 0.007026598725856425), (3, 0.020400080447157378), (6, 0.003728010614697474), (8, 0.00948386191556759), (11, 0.025343293191890992), (6, 0.006344835080904147), (11, 0.0159827220965609), (14, 0.011144264813696745), (15, 0.007331014842165503), (28, 0.0014497153882413207), (5, 0.021969432881049394), (6, 0.008911633378357455), (7, 0.014521089470083912), (9, 0.024701005937586898), (11, 0.01282725939651922), (5, 0.046116965595231596), (6, 0.00267239806801576), (9, 0.012962731527460972), (10, 0.011538190997504824), (11, 0.026927221973684842), (1, 0.05724451153936908), (5, 0.016877938375748296), (6, 0.004890236090479306), (8, 0.00737568554728487), (14, 0.006871489682412035), (0, 0.003514834230373475), (6, 0.001930529190878152), (14, 0.006781682050340082), (17, 0.006302307729400447), (18, 0.0575034236121058), (0, 0.004046646823944413), (7, 0.012675840393931688), (11, 0.022395317784263406), (14, 0.007807786135142372), (15, 0.005136184126819541), (0, 0.004345810285563815), (2, 0.013336589713368138), (6, 0.002386944323518779), (11, 0.01202548759622852), (14, 0.00838500585047696), (2, 0.01282168326189491), (5, 0.019800316013805338), (9, 0.02226219156644913), (10, 0.029723529077097657), (14, 0.01612254579477476), (3, 0.016924440534952175), (5, 0.01800464931537239), (6, 0.0010433381603640788), (10, 0.018018625463540133), (14, 0.02932802325500092), (7, 0.0001000001120171377) (3, 0.01310000777161024) (5, 0.010700000000000000) (6, 0.001000000000000000) (8, 0.001000000000000000) (10, 0.001000000000000000) (12, 0.001000000000000000) (14, 0.001000000000000000) (16, 0.001000000000000000) (18, 0.001000000000000000) (20, 0.001000000000000000) (22, 0.001000000000000000) (24, 0.001000000000000000) (26, 0.001000000000000000) (28, 0.001000000000000000) (30, 0.001000000000000000) (32, 0.001000000000000000) (34, 0.001000000000000000) (36, 0.001000000000000000) (38, 0.001000000000000000) (40, 0.001000000000000000) (42, 0.001000000000000000) (44, 0.001000000000000000) (46, 0.001000000000000000) (48, 0.001000000000000000) (50, 0.001000000000000000) (52, 0.001000000000000000) (54, 0.001000000000000000) (56, 0.001000000000000000) (58, 0.001000000000000000) (60, 0.001000000000000000) (62, 0.001000000000000000) (64, 0.001000000000000000) (66, 0.001000000000000000) (68, 0.001000000000000000) (70, 0.001000000000000000) (72, 0.001000000000000000) (74, 0.001000000000000000) (76, 0.001000000000000000) (78, 0.001000000000000000) (80, 0.001000000000000000) (82, 0.001000000000000000) (84, 0.001000000000000000) (86, 0.001000000000000000) (88, 0.001000000000000000) (90, 0.001000000000000000) (92, 0.001000000000000000) (94, 0.001000000000000000) (96, 0.001000000000000000) (98, 0.001000000000000000) (100, 0.001000000000000000) (102, 0.001000000000000000) (104, 0.001000000000000000) (106, 0.001000000000000000) (108, 0.001000000000000000) (110, 0.001000000000000000) (112, 0.001000000000000000) (114, 0.001000000000000000) (116, 0.001000000000000000) (118, 0.001000000000000000) (120, 0.001000000000000000) (122, 0.001000000000000000) (124, 0.001000000000000000) (126, 0.001000000000000000) (128, 0.001000000000000000) (130, 0.001000000000000000) (132, 0.001000000000000000) (134, 0.001000000000000000) (136, 0.001000000000000000) (138, 0.001000000000000000) (140, 0.001000000000000000) (142, 0.001000000000000000) (144, 0.001000000000000000) (146, 0.001000000000000000) (148, 0.001000000000000000) (150, 0.001000000000000000) (152, 0.001000000000000000) (154, 0.001000000000000000) (156, 0.001000000000000000) (158, 0.001000000000000000) (160, 0.001000000000000000) (162, 0.001000000000000000) (164, 0.001000000000000000) (166, 0.001000000000000000) (168, 0.001000000000000000) (170, 0.001000000000000000) (172, 0.001000000000000000) (174, 0.001000000000000000) (176, 0.001000000000000000) (178, 0.001000000000000000) (180, 0.001000000000000000) (182, 0.001000000000000000) (184, 0.001000000000000000) (186, 0.001000000000000000) (188, 0.001000000000000000) (190, 0.001000000000000000) (192, 0.001000000000000000) (194, 0.001000000000000000) (196, 0.001000000000000000) (198, 0.001000000000000000) (200, 0.001000000000000000) (202, 0.001000000000000000) (204, 0.001000000000000000) (206, 0.001000000000000000) (208, 0.001000000000000000) (210, 0.001000000000000000) (212, 0.001000000000000000) (214, 0.001000000000000000) (216, 0.001000000000000000) (218, 0.001000000000000000) (220, 0.001000000000000000) (222, 0.001000000000000000) (224, 0.001000000000000000) (226, 0.001000000000000000) (228, 0.001000000000000000) (230, 0.001000000000000000) (232, 0.001000000000000000) (234, 0.001000000000000000) (236, 0.001000000000000000) (238, 0.001000000000000000) (240, 0.001000000000000000) (242, 0.001000000000000000) (244, 0.001000000000000000) (246, 0.001000000000000000) (248, 0.001000000000000000) (250, 0.001000000000000000) (252, 0.001000000000000000) (254, 0.001000000000000000) (256, 0.001000000000000000) (258, 0.001000000000000000) (260, 0.001000000000000000) (262, 0.001000000000000000) (264, 0.001000000000000000) (266, 0.001000000000000000) (268, 0.001000000000000000) (270, 0.001000000000000000) (272, 0.001000000000000000) (274, 0.001000000000000000) (276, 0.001000000000000000) (278, 0.001000000000000000) (280, 0.001000000000000000) (282, 0.001000000000000000) (284, 0.001000000000000000) (286, 0.001000000000000000) (288, 0.001000000000000000) (290, 0.001000000000000000) (292, 0.001000000000000000) (294, 0.001000000000000000) (296, 0.001000000000000000) (298, 0.001000000000000000) (300, 0.001000000000000000) (302, 0.001000000000000000) (304, 0.001000000000000000) (306, 0.001000000000000000) (308, 0.001000000000000000) (310, 0.001000000000000000) (312, 0.001000000000000000) (314, 0.001000000000000000) (316, 0.001000000000000000) (318, 0.001000000000000000) (320, 0.001000000000000000) (322, 0.001000000000000000) (324, 0.001000000000000000) (326, 0.001000000000000000) (328, 0.001000000000000000) (330, 0.001000000000000000) (332, 0.001000000000000000) (334, 0.001000000000000000) (336, 0.001000000000000000) (338, 0.001000000000000000) (340, 0.001000000000000000) (342, 0.001000000000000000) (344, 0.001000000000000000) (346, 0.001000000000000000) (348, 0.001000000000000000) (350, 0.001000000000000000) (352, 0.001000000000000000) (354, 0.001000000000000000) (356, 0.001000000000000000) (358, 0.001000000000000000) (360, 0.001000000000000000) (362, 0.001000000000000000) (364, 0.001000000000000000) (366, 0.001000000000000000) (368, 0.001000000000000000) (370, 0.001000000000000000) (372, 0.001000000000000000) (374, 0.001000000000000000) (376, 0.001000000000000000) (378, 0.001000000000000000) (380, 0.001000000000000000) (382, 0.001000000000000000) (384, 0.001000000000000000) (386, 0.001000000000000000) (388, 0.001000000000000000) (390, 0.001000000000000000) (392, 0.001000000000000000) (394, 0.001000000000000000) (396, 0.001000000000000000) (398, 0.001000000000000000) (400, 0.001000000000000000) (402, 0.001000000000000000) (404, 0.001000000000000000) (406, 0.001000000000000000) (408, 0.001000000000000000) (410, 0.001000000000000000) (412, 0.001000000000000000) (414, 0.001000000000000000) (416, 0.001000000000000000) (418, 0.001000000000000000) (420, 0.001000000000000000) (422, 0.001000000000000000) (424, 0.001000000000000000) (426, 0.001000000000000000) (428, 0.001000000000000000) (430, 0.001000000000000000) (432, 0.001000000000000000) (434, 0.001000000000000000) (436, 0.001000000000000000) (438, 0.001000000000000000) (440, 0.001000000000000000) (442, 0.001000000000000000) (444, 0.001000000000000000) (446, 0.001000000000000000) (448, 0.001000000000000000) (450, 0.001000000000000000) (452, 0.001000000000000000) (454, 0.001000000000000000) (456, 0.001000000000000000) (458, 0.001000000000000000) (460, 0.001000000000000000) (462, 0.001000000000000000) (464, 0.001000000000000000) (466, 0.001000000000000000) (468, 0.001000000000000000) (470, 0.001000000000000000) (472, 0.001000000000000000) (474, 0.001000000000000000) (476, 0.001000000000000000) (478, 0.001000000000000000) (480, 0.001000000000000000) (482, 0.001000000000000000) (484, 0.001000000000000000) (486, 0.001000000000000000) (488, 0.001000000000000000) (490, 0.001000000000000000) (492, 0.001000000000000000) (494, 0.001000000000000000) (496, 0.001000000000000000) (498, 0.001000000000000000) (500, 0.001000000000000000) (502, 0.001000000000000000) (504, 0.001000000000000000) (506, 0.001000000000000000) (508, 0.001000000000000000) (510, 0.001000000000000000) (512, 0.001000000000000000) (514, 0.001000000000000000) (516, 0.001000000000000000) (518, 0.001000000000000000) (520, 0.001000000000000000) (522, 0.001000000000000000) (524, 0.001000000000000000) (526, 0.001000000000000000) (528, 0.001000000000000000) (530, 0.001000000000000000) (532, 0.001000000000000000) (534, 0.001000000000000000) (536, 0.001000000000000000) (538, 0.001000000000000000) (540, 0.001000000000000000) (542, 0.001000000000000000) (544, 0.001000000000000000) (546, 0.001000000000000000) (548, 0.001000000000000000) (550, 0.001000000000000000) (552, 0.001000000000000000) (554, 0.001000000000000000) (556, 0.001000000000000000) (558, 0.001000000000000000) (560, 0.001000000000000000) (562, 0.001000000000000000) (564, 0.001000000000000000) (566, 0.001000000000000000) (568, 0.001000000000000000) (570, 0.001000000000000000) (572, 0.001000000000000000) (574, 0.001000000000000000) (576, 0.001000000000000000) (578, 0.001000000000000000) (580, 0.001000000000000000) (582, 0.001000000000000000) (584, 0.001000000000000000) (586, 0.001000000000000000) (588, 0.001000000000000000) (590, 0.001000000000000000) (592, 0.001000000000000000) (594, 0.001000000000000000) (596, 0.001000000000000000) (598, 0.001000000000000000) (600, 0.001000000000000000) (602, 0.001000000000000000) (604, 0.001000000000000000) (606, 0.001000000000000000) (608, 0.001000000000000000) (610, 0.001000000000000000) (612, 0.001000000000000000) (614, 0.001000000000000000) (616, 0.001000000000000000) (618, 0.001000000000000000) (620, 0.001000000000000000) (622, 0.001000000000000000) (624, 0.001000000000000000) (626, 0.001000000000000000) (628, 0.001000000000000000) (630, 0.001000000000000000) (632, 0.001000000000000000) (634, 0.001000000000000000) (636, 0.001000000000000000) (638, 0.001000000000000000) (640, 0.001000000000000000) (642, 0.001000000000000000) (644, 0.001000000000000000) (646, 0.001000000000000000) (648, 0.001000000000000000) (650, 0.001000000000000000) (652, 0.0010
```

following 1817, Mexico and the Gulf was highly used due to tension which ultimately lead to treaty between the two nation.

2. The second topic has keywords Jobs, budget,economic, program etc. clearly depicts the economic reforms, programs to increase production and Job in the state. This topic was dominant around 1850s.
3. The second topic and third topic does not differ too much.
4. The fourth topic has varied keywords, therefore, cannot be annotated. It does not give details about any significant topic except economy like cents, gold.farms etc.
5. Fifth topic has keywords has interstate, silver, iraq, terrorist, soviet etc. As we know, United States of America has almost 50 states, In beginning of its structural development, there were many disputes between states and many reforms and rules were made to make peace and for the betterment of all the states.
6. Sixth topic had keywords iraq,al,terrorist,saddam etc. This obviously is the major era of the United States history. This is around 2003 when US invaded iraq. It had also involvement of Afghanistan.
7. Seventh topic comprises of keywords texas, kansas, constitution ,spain ,slavery. This topic is very much generalised which talk about basic achievements and progress states and the cities have made. Constitution word is used frequently because any amendment in the constitution has a say in union address.
8. Topic Number 8 has keywords bank, vietnam ,veterans ,banks etc. This is basically around 1974 when the vietnam war happened. If we go through csv file vietnam came into addresses in the year 1966 and continued upto 1990s. In one of the year whole depiction of the war is given.
9. Ninth topic includes notes, coinage, gold,silver,currency,notes. These words are very much frequent because it comprises of the basic financial system of the country.
10. Tenth topic has keyword enemy,japanese ,oil etc. This topic does not get focussed to any specific point but gives the hint of tension between Japan and USA before world war 2.

3 LDA topic modeling

Above exercise is repeated using Latent Dirichlet Allocation on the data. Overall it is found that LDA is better for visualization than LSI. Basically it covers more important topics than lsi model has covered. All the observations are mentioned in the following points. Let us consider first ten topics in the LDA model also.

1. First topic has keywords iraq,enemy,terrorists etc. This is in resemblance with the sixth topic of lsi model.

```

[[0,
'0.000**iraq' + 0.000**tonight' + 0.000**terrorists' + 0.000**jobs' + 0.000**americans' + 0.000**program' + 0.000**america' + 0.000**workers' + 0.000**terror' + 0.000**help"),
(1,
'0.000**tonight' + 0.000**jobs' + 0.000**americans' + 0.000**inflation' + 0.000**help' + 0.000**child' + 0.000**college' + 0.000**economic' + 0.000**children' + 0.000**programs"),
(2,
'0.000**mexico' + 0.000**award' + 0.000**treaty' + 0.000**majesty' + 0.000**mexican' + 0.000**texas' + 0.000**june' + 0.000**upon' + 0.000**freedmen' + 0.000**cent"),
(3,
'0.000**soviet' + 0.000**economic' + 0.000**jobs' + 0.000**texas' + 0.000**today' + 0.000**upon' + 0.000**subject' + 0.000**americans' + 0.000**cable' + 0.000**programs"),
(4,
'0.000**banks' + 0.000**tonight' + 0.000**subject' + 0.000**silver' + 0.000**upon' + 0.000**treasury' + 0.000**challenge' + 0.000**militia' + 0.000**americans' + 0.000**bank"),
(5,
'0.000**vietnam' + 0.000**spain' + 0.000**billion' + 0.000**article' + 0.000**programs' + 0.000**gentlemen' + 0.000**vessels' + 0.000**percent' + 0.000**program' + 0.000**economic"),
(6,
'0.000**soviet' + 0.000**oil' + 0.000**program' + 0.000**budget' + 0.000**cent' + 0.000**afghanistan' + 0.000**help' + 0.000**emancipation' + 0.000**tons' + 0.000**programs"),
(7,
'0.000**gentlemen' + 0.000**budget' + 0.000**upon' + 0.000**today' + 0.000**program' + 0.000**help' + 0.000**let' + 0.000**tonight' + 0.000**economic' + 0.000**oil"),
(8,
'0.000**economic' + 0.000**kansas' + 0.000**banks' + 0.000**democracy' + 0.000**cent' + 0.000**veterans' + 0.000**upon' + 0.000**atomic' + 0.000**communist' + 0.000**per"),
(9,
'0.000**program' + 0.000**programs' + 0.000**billion' + 0.000**gold' + 0.000**economic' + 0.000**tonight' + 0.000**budget' + 0.000**help' + 0.000**federal' + 0.000**silver"),
(10,
'0.000**tonight' + 0.000**economic' + 0.000**americans' + 0.000**budget' + 0.000**jobs' + 0.000**help' + 0.000**upon' + 0.000**recovery' + 0.000**attentions' + 0.000**militia"),
(11,
'0.000**tonight' + 0.000**americans' + 0.000**democracy' + 0.000**help' + 0.000**jobs' + 0.000**programs' + 0.000**billion' + 0.000**program' + 0.000**cannot' + 0.000**today"),
(12,
'0.001**tonight' + 0.000**vietnam' + 0.000**mexico' + 0.000**billion' + 0.000**budget' + 0.000**help' + 0.000**gold' + 0.000**today' + 0.000**program' + 0.000**upon"),
(13,
'0.000**spain' + 0.000**gold' + 0.000**kansas' + 0.000**treaty' + 0.000**currency' + 0.000**paper' + 0.000**upon' + 0.000**constitution' + 0.000**silver' + 0.000**minister"),
(14,
'0.000**economic' + 0.000**hitler' + 0.000**program' + 0.000**gentlemen' + 0.000**upon' + 0.000**texas' + 0.000**mexico' + 0.000**interstate' + 0.000**statute' + 0.000**industrial")]

```

Figure 3: LDA-Topics

2. Second topic has keywords Jobs, Americans, inflation, economic ,college. This topic has resemblance with second topic of lsi model but it can be clearly observed that words are well interconnected and covers many of the topics in which reforms have been made.
3. Third topic consist of mexico, treaty ,award, mexican etc. This directly gives reference to the Treaty of Guadalupe Hidalgo, signed on February 2, 1848, ended the Mexican-American War in favor of the United States.
4. Fourth topic has varied words which do not depict a major event in the history.
5. The fifth topic has keywords bank, treasure,bank ,militia etc. This topic hints about the financial system.
6. The sixth topic has keywords vietnam, spain, billions,articles ,programs. This also different words which means the speech did not focus on a particular event.
7. Seventh topic consist of keywords soviet ,oil ,budget, cent ,aghanistan. Basically soviet union was formed in 1922 and it is one of the major events in world history and therefore it has much frequency in the Union addresses.
8. Topic number 8 has common keywords. This topic is dominant in those years in which focuses on general topics like in initial years.
9. Ninth topic of lda model resembles with the seventh topic. It has keywords economic, banks ,democracy,cent etc. Kansas city is an important city in the USA due to its aviation manufacturing and revenue generation, Moreover, there is a high amount of mention of Kansas Pacific in 1900s due to some financial debts in the area.
10. The 14th topic has keywords hitler,economic,texas,mexico. This is the time of World War 1 when Hitler and its Nazi germans ruled in some parts of Europe. This topic is not mentioned with significant amount of weight in lsi. Therefore, lda model proves to be better in the data of state of the union.csv



Figure 4: Some important visualizations

4 Change in Speeches overtime

There are many ways to summarize the changes overtime. this. Possibilities include: visualizations, grouping speeches by decade after topic modeling, and grouping speeches by decade before topic modeling.LDA model has been taken as a base for the algorithm. pyLDAvis ,an inbuilt pyton package, which is commonly used in topic modeling visualization proves to be very helpful in this case. This is the best tool according to many of the data scientists. Left-hand side contains bubbles which represent topics. A good topic will have big, non-overlapping bubble scattered throughout the chart. If we move the mouse pointer over one of the bubbles, the words order will change accordingly and bars will be updated and some part of it will be red showing percentage of word in the content . These words are the salient keywords that form the selected topic. The corresponding visualization is shown in **Figure:4**. Using pyLDAvis notebook and topics of LDA following observations can be recognised-

- US got independent and 1770s and in years following it, speeches include more about reforms and programs etc.

- On February 2, 1848, the Treaty of Guadalupe Hidalgo was signed between Mexico and USA which brought to a close the Mexican-American War in 1850s when Mr. James Polk was the president. It can be easily observed mexican,mexicans and treaty were important words in 1850s.
- In 1900s President Theodore Roosevelt had many progressive goals in terms of economy,financial and political power. Therefore, the topics include interstate, states, bank ,cent,notes etc.
- In 1910s, World War-I raged in Europe, but the President avoided any involvement at first but then were part of allies. Topic which includes war,enemy,hitler, europe were dominant in this era.
- In 1920s came the year of great depression which lead to stock market crash. In the speeches of these years, topic mainly consisted of budget, economic,reforms , jobs, programs,gold and words related to finance.
- In the era between great depression and world-war-II, many reforms were made in fields related to agriculture, industry, finance etc. by the president Franklin D. Roosevelt.
- 1940s was the era of world-war-II. In these years speeches had too much content related to japanese, enemy,assassination,killing,war, atomic("Atomic bomb"). It was clear from the content, it was all war in these times.
- In the year 1950s and 1960s, topic which had keywords militia,tactics,vietnam,war were dominant because of Vietnam war which also lead to civil rights movement.
- The watergate scandal cannot be determined in lda model. Topics related to it cannot get recognised through this algorithm.
- In 2000s, there were two major events one was 9/11 attack on the world trade centre. Topics related to terror, killing,terrorists were common in these years.
- In 2003, US invaded iraq and therefore there are many topics which had iraq, iraqi ,saddam with significant amount of weight. The President at that time was George W. Bush. Therefore, many historical events in US can easily be recognised through this model.

5 Different Document Set

LDA is applied on different document set which is a collection of AP wire stories. Same process is repeated as for the state of the union dataset and results of both are compared. If we go through the csv file, we can easily observe that they are news articles. Now the news has widespread topics, therefore,hundred number of topics were chosen. Lets here also consider first 10 topics of lda model.

1. Now, the first topic has main keywords sales, taxes,billions,india, delhi. Sales and taxes are obviously will have more weights they are a core part of any nation. India has come into consideration as after going through the data, it is clearly visible that there is too much mention of indians slavery and indians residing in the nation. Moreover war between India and Pakistan is discussed.

```

[[41,
'0.014**sales" + 0.009**tax" + 0.009**agreement" + 0.007**united" + 0.006**abuse" + 0.006**billion" + 0.005**delhi" + 0.005**news" + 0.005**india" + 0.005**joints"'),
(82,
'0.014**primary" + 0.010**cars" + 0.008**plant" + 0.008**rocks" + 0.007**domestic" + 0.007**jefferson" + 0.007**williams" + 0.006**techniques" + 0.005**bush" + 0.005**pont"'),
(38,
'0.013**sasser" + 0.013**climbed" + 0.011**barr" + 0.010**cosby" + 0.008**november" + 0.007**adopt" + 0.006**strains" + 0.006**afghanistan" + 0.006**welcome" + 0.006**realistic"'),
(63,
'0.029**presley" + 0.017**disc" + 0.016**happens" + 0.013**elvis" + 0.011**le" + 0.010**french" + 0.007**love" + 0.007**extreme" + 0.007**article" + 0.006**sank"'),
(33,
'0.013**interstate" + 0.009**epa" + 0.008**miller" + 0.007**hong" + 0.007**kong" + 0.007**tariffs" + 0.007**chinese" + 0.006**connection" + 0.006**arrested" + 0.006**earned"'),
(49,
'0.009**economy" + 0.008**study" + 0.006**hoffman" + 0.006**percent" + 0.005**houston" + 0.005**groups" + 0.005**agencies" + 0.005**part" + 0.005**toronto" + 0.005**usually"'),
(94,
'0.011**energy" + 0.007**oh" + 0.007**dress" + 0.006**women" + 0.006**committee" + 0.006**educated" + 0.005**ourselves" + 0.005**salinas" + 0.005**chernobyl" + 0.005**she"'),
(10,
'0.016**thompson" + 0.007**wilson" + 0.007**delegates" + 0.006**now" + 0.006**ala" + 0.006**air" + 0.005**iran" + 0.005**kashmir" + 0.005**speed" + 0.004**literary"'),
(90,
'0.007**below" + 0.006**she" + 0.006**deaths" + 0.006**tickets" + 0.005**texas" + 0.005**ms" + 0.005**there" + 0.004**state" + 0.004**night" + 0.004**mexico"'),
(83,
'0.018**video" + 0.015**speaker" + 0.012**indictment" + 0.010**ohio" + 0.010**cocaine" + 0.008**sing" + 0.007**smuggler" + 0.007**ethics" + 0.007**committee" + 0.006**wage"'),
(14,
'0.006**panama" + 0.006**lafayette" + 0.006**late" + 0.006**gold" + 0.006**down" + 0.006**together" + 0.006**london" + 0.005**noriega" + 0.005**technology" + 0.005**square"'),
(69,
'0.047**ec" + 0.018**bendjedid" + 0.015**hamilton" + 0.012**buses" + 0.010**sweden" + 0.008**blast" + 0.007**deductions" + 0.007**nelson" + 0.006**legislature" + 0.005**drivers"'),
(48,
'0.012**awards" + 0.011**kohl" + 0.008**flood" + 0.008**deer" + 0.007**retire" + 0.007**white" + 0.007**chancellor" + 0.006**waters" + 0.006**intervention" + 0.005**random"'),
(39,
'0.017**trash" + 0.014**mecham" + 0.013**cent" + 0.013**roberts" + 0.010**bushel" + 0.008**lower" + 0.007**victims" + 0.007**higher" + 0.006**name" + 0.006**futures"'),
(98,
'0.014**navy" + 0.012**ships" + 0.011**sea" + 0.009**computer" + 0.009**detained" + 0.007**sony" + 0.007**expenses" + 0.006**tools" + 0.006**rehabilitation" + 0.006**straight"'),
(6,
'0.012**shearson" + 0.008**draws" + 0.007**sweden" + 0.007**starring" + 0.007**allen" + 0.006**winners" + 0.006**producer" + 0.006**undecided" + 0.005**jonathan" + 0.005**gun"')],

```

Figure 5: Topics of the second document

- The second topic has keywords plant, cars and some names. After going through the data, there many news which involve accidents, car chase of police and train accidents. These events are described thoroughly and therefore, has a high weight. Names are basically of presidents of USA. News always mention the name of the President when major event happens.
- The third topic also consists of a name sasser. He was United States Senator. But he other words are common english words which are commonly used. Hence, this topic was not useful.
- Fourth topic consists of keywords presley, disc, happens, elvis, french, love. This topic has varied words but some can be connected like Elvis Presley is an American singer whose life is too many discussed. The words like love and french are frequently used.
- Fifth topic can be annotated as interstate, miller, hong kong, tariffs, chinese etc. Miller is again a famous Personality. Hong Kong and China have been in news for many of the reasons.
- Sixth topic can be annotated as economy study houston groups agencies. This topic mainly covers the news which are related economy and its future.
- Seventh Topic has keywords women, dress, energy, chernobyl, she. Whenever topic related to women dresses and education are discussed this topic is dominant. Moreover, it has word chernobyl which is related to the disaster in chernobyl and its consequences around that region.
- Eighth topic is a bit wider one as it comprises of Thomson, Wilson, delegates, kashmir, iran. First two are again names. Iran and US has a long history which are frequently described in the news. Basically, Iranian fighter planes, killings, islamic revolution were the main reasons due to which a few news are based on Iran.

9. Ninth topic again does not depict any significant news or event. It consists of common english words.
10. Tenth Topic has important keywords like video, smuggler, Ohio, cocaine, indictment, speaker. This topic is clearly about drug smuggling and serious punishment which was given for the crime. Videos were part of evidence in only a few of the cases. The speaker word is mainly due to speaker of the house which is an important person in the country.

Now, to compare the two lda model. Its obvious one news article focuses on a single topic while union speech covers all important reforms and events occurred. Therefore, in the second data, number of topics has to be increased and if we look at the dominant topics of each row, easily we can two or more words that can be the topic. But in union addresses some topics, just does not represent the content of the speech. Some topics even had no clear concept. Increasing the number of topics also did not help because of varied content in the speech. Overall, due to reasons mentioned, **output of AP wire collection document proves to be the better one.**