

LING 410X: Language as Data

Semester: Spring '18

Instructor: Sowmya Vajjala

Iowa State University, USA

27 March 2018

Class Outline - Topic Modeling in detail

- ▶ Continuing from last class: how does a topic model "learn"?
- ▶ Sensitivity of topic model to its parameters (e.g., num. iterations)
- ▶ Evaluating topic models automatically
- ▶ Human evaluation of topic models

Class Outline - Topic Modeling in detail

- ▶ Continuing from last class: how does a topic model "learn"?
- ▶ Sensitivity of topic model to its parameters (e.g., num. iterations)
- ▶ Evaluating topic models automatically
- ▶ Human evaluation of topic models
- ▶ Reminder: A5 submission due on 3/31. If topicmodels is slow, or you don't want to learn an extra library, you can also do the assignment with `mallet` instead.

Exercise on Thursday

- ▶ How many of tried this out?
- ▶ How many of you successfully make this work on your computers?

Exercise on Thursday

- ▶ How many of tried this out?
- ▶ How many of you successfully make this work on your computers?
- ▶ The task was to "learn" a topic model, based on a dataset.

"Learning" a topic model: intuitive explanation of LDA

- ▶ Aim: organize a topic collection into topics, where each topic is a collection of words.
- ▶ Assumptions: each document is a mixture of topics, each topic is a mixture of keywords, a keyword can exist in multiple topics.

"Learning" a topic model: intuitive explanation of LDA

- ▶ Aim: organize a topic collection into topics, where each topic is a collection of words.
- ▶ Assumptions: each document is a mixture of topics, each topic is a mixture of keywords, a keyword can exist in multiple topics.
- ▶ Learning:
 - ▶ Decide on the number of topics K .
 - ▶ Go through each document, randomly assign each word in a document to one of the K topics.
 - ▶ At that point, you have on poorly represented topic model already, where each document is represented as a collection of topics.

after random initialization

- ▶ Again, we start revisiting all documents.
- ▶ For each document d , for each word w in it, we compute two probabilities:

$P(t|d)$ i.e, proportion of words in the document d assigned topic t

$P(w|t)$ i.e., proportion of documents with topic t , containing word w .

after random initialization

- ▶ Again, we start revisiting all documents.
- ▶ For each document d , for each word w in it, we compute two probabilities:
 - $P(t|d)$ i.e, proportion of words in the document d assigned topic t
 - $P(w|t)$ i.e., proportion of documents with topic t , containing word w .
- ▶ We reassign w a new topic t with a probability $P(t|d) * P(w|t)$
- ▶ Assumption made: each time we assign a word a new topic like this, we are assuming all other assignments are correct.

after random initialization

- ▶ Again, we start revisiting all documents.
- ▶ For each document d , for each word w in it, we compute two probabilities:
 - $P(t|d)$ i.e, proportion of words in the document d assigned topic t
 - $P(w|t)$ i.e., proportion of documents with topic t , containing word w .
- ▶ We reassign w a new topic t with a probability $P(t|d) * P(w|t)$
- ▶ Assumption made: each time we assign a word a new topic like this, we are assuming all other assignments are correct.

Finally,

- ▶ If you keep doing this process several times (iterations), you will reach a state where there is not much change in the topic-word distributions between iterations.
- ▶ At that point, one final time, you go through all documents and assign topic probabilities to them.

Note: Read <https://goo.gl/jbgCKq> for more detailed explanation of this.

Some name dropping

- ▶ The topics and words are assumed to have "Dirichlet distribution" which is one of the several probability distributions studied in probability and statistics (D of LDA comes from there!)
- ▶ The iterative process I conceptually explained is generally implemented by an algorithm called "Gibbs Sampling".

"Learning" a topic model: steps in the code

1. Pre-processing of the corpus (in author's version: splitting it into chunks, pre-processing of chunks)
2. Getting it into a two column format (id, text).
3. Using `mallet` package in R to build a topic model.
4. use `mallet.import()` function to convert a dataset into `mallet` format
5. use `MalletLDA()` function to create a topic model, setting its parameters.
6. observe and analyze the output in different ways.

(Let me go through this taking our movie review data from text classification weeks)

Step 1: Pre-processing of the corpus

```
setwd("~/Dropbox/ClassroomSlides-BothCourses/LING410X/21Mar2018/data/reviews")

get_text_string <- function(file_path)
{
  fulltext <- scan(file_path, what = "character", sep = "\n")
  fulltext_as_string <- tolower(paste(fulltext, collapse = " "))
  return (fulltext_as_string)
}

files.v <- dir(path=getwd(), pattern="*.txt")
documents = c()
for(i in 1:length(files.v)){
  documents = c(documents, get_text_string(files.v[i]))
}
```

Step 2: Getting it into a two column format (id, text)

```
docids = seq(1:length(documents))  
  
fortopicmodels <- cbind(docids,documents)  
finaldocuments <- as.data.frame(fortopicmodels, stringsAsFactors=F)
```

Step 3: Import data into mallet

```
stoplistfile = "/ClassroomSlides-BothCourses/LING410X/21Mar2018/stoplist.csv"
mallet.instances <- mallet.import(finaldocuments$docids,
                                  finaldocuments$documents,
                                  stoplistfile,FALSE,token.regexp="[\\p{L}']+")
```

- ▶ What does that regular expression mean?

Step 3: Import data into mallet

```
stoplistfile = "/ClassroomSlides-BothCourses/LING410X/21Mar2018/stoplist.csv"
mallet.instances <- mallet.import(finaldocuments$docids,
                                  finaldocuments$documents,
                                  stoplistfile,FALSE,token.regexp="[\\p{L}']+")
```

- ▶ What does that regular expression mean? It means: match one or more occurrences of any Unicode character or a " ' " .
- ▶ What is that "FALSE"?

Step 3: Import data into mallet

```
stoplistfile = "/ClassroomSlides-BothCourses/LING410X/21Mar2018/stoplist.csv"
mallet.instances <- mallet.import(finaldocuments$docids,
                                  finaldocuments$documents,
                                  stoplistfile,FALSE,token.regexp="[\\p{L}']+")
```

- ▶ What does that regular expression mean? It means: match one or more occurrences of any Unicode character or a " ' " .
- ▶ What is that "FALSE"? - It is the argument "preserve case".

More on regular expressions:

<http://www.regular-expressions.info>

More on unicode in regular expressions:

<http://www.regular-expressions.info/unicode.html>

Setting up the topic model training process

```
topic.model <- MalletLDA(num.topics=XX) #You have to specify this.  
topic.model$loadDocuments(mallet.instances) #from previous slide.  
topic.model$setAlphaOptimization(40, 80) #Optional  
topic.model$train(NUM) #Actual training happens here.
```

What are all those parameters?

- ▶ Best place to look for that= ?MalletLDA
- ▶ `setAlphaOptimization(40,80)` means: perform model optimization after every 40 iterations, starting first after 80 iterations.
- ▶ number of iterations: the more the better, generally. However, trade off is that training becomes very slow.
- ▶ Every 50 iterations, R prints a summary of your topic model showing top 7 words per topic.
- ▶ Every 10 iterations, R also shows a log likelihood of this topic model (the closer to 0, the better).

Output screenshots for Thursday's Data, when I
trained it last year

Output screenshot - at the beginning

```
Console - 1
> #This starts the training process
> topic.model.train(400)
Mar 28, 2017 8:55:31 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <1> LL/token: -10.34105
Mar 28, 2017 8:55:33 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <2> LL/token: -9.85248
Mar 28, 2017 8:55:35 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <3> LL/token: -9.6863
Mar 28, 2017 8:55:36 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <4> LL/token: -9.59737
Mar 28, 2017 8:55:38 AM cc.mallet.topics.ParallelTopicModel estimate
INFO:
0      0.11628 doctor made thought great mind letter long
1      0.11628 irish ireland mr family native father called
2      0.11628 sir reilly man robert sauire folliard mr
3      0.11628 d'arcy bejaobers camp gold demond man se
4      0.11628 rachael warren magpie miss george gregory woman
5      0.11628 abellina time men thou long frazier ground
6      0.11628 lady young read sort fair table mr
7      0.11628 ' i miss 'and wylder 'you rachel
8      0.11628 property man justice church principles fact judge
9      0.11628 don parker farrel horse pablo mike miguel
10     0.11628 heart love father mind tears girl child
11     0.11628 father heart mother god cannon son man
12     0.11628 bryce cardigan colonel shirley pennington sequoia woods
13     0.11628 god good poor day father make home
14     0.11628 nall tie sheila danny children prairie wild
15     0.11628 door room night hand madame window house
16     0.11628 men lord people country cumber friend good
17     0.11628 ye mr story w' eagle gilbert henry
18     0.11628 peter mary dinner day great emily made
19     0.11628 men people country state land indians english
20     0.11628 day young life children years death family
21     0.11628 sir mr replied good hycy man bryan
22     0.11628 night face life arms clarence found long
23     0.11628 man country time house fact blood found
24     0.11628 uncle lady could milly time mad good
25     0.11628 day long work office hand asked home
26     0.11628 high river small large green trees years
27     0.11628 charles maria french washington great place make
28     0.11628 mr sir captain lake man larkin mark
29     0.11628 wid good man wuld none priest night
30     0.11628 night truth matter felt poor replied make
31     0.11628 big mother flurry water house family back
32     0.11628 replied poor god man exclaimed ay good
33     0.11628 room eyes morston mrs face stood passed
34     0.11628 john money good dollars years pay day
35     0.11628 susan mortie mrs billy sally asked miss
36     0.11628 doctor sir mrs toole puddock sturk dangerfield
37     0.11628 gerald san city years race young francisco
38     0.11628 rather young margaret lydia good aunt mrs
39     0.11628 replied miss harry woodard family mr mother
40     0.11628 night made men man good knew life
41     0.11628 spirit denis light father deep spoke appearance
42     0.11628 school van time man boys job horses

Mar 28, 2017 8:55:38 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <5> LL/token: -9.54878
Mar 28, 2017 8:55:39 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <6> LL/token: -9.50984
```

Output screenshot - after 400 iterations

```
INFO: <368> LL/token: -9.2627
Mar 28, 2017 8:56:30 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <378> LL/token: -9.25952
Mar 28, 2017 8:56:32 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <388> LL/token: -9.25652
Mar 28, 2017 8:56:33 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <398> LL/token: -9.25065
Mar 28, 2017 8:56:35 AM cc.mallet.topics.ParallelTopicModel estimate
INFO:
0      0.2753 doctor letter time read great thought long
1      0.01875 irish ireland family mr john county race
2      0.02243 sir reilly robert mr folliard replied whitecraft
3      0.01584 d'arcy bejaubers don camp demod gold jos
4      0.01545 rachael warren magpie billy clarence gregory alice
5      0.01169 abellino rasabella flodoardo andreas venice thou dage
6      0.3172 mr miss young good lady man great
7      0.11139 ' i 'and poor 'you 'well looked
8      0.10621 judge gentlemen low court man justice protestant
9      0.01702 don parker farrel pablo mike miguel key
10     0.29702 heart love father god mother tears girl
11     0.06827 connor son mother fardoragha una heart battle
12     0.01467 bryce cardigan shirley colonel pennington timber sequoia
13     0.06262 father priest church denis bishop poor holy
14     0.01356 nell tin sheila danny children prairie johann
15     0.32824 door room night house window bed time
16     0.07018 sir mr phil darby lord n'clutchy val
17     0.01095 ye wi' o' hae na gilbert sae
18     0.05661 susan miss peter mary sue emily ella
19     0.0093 nen irish san california state city land
20     0.34224 life nan country people young day years
21     0.02738 hycy bryan m'mohan burke kathleen replied cavanagh
22     0.3506 eyes face night life voice heart back
23     0.16512 country nen house purcel night people party
24     0.0184 radame uncle milly moud cousin silas mary
25     0.38849 day time home back long told make
26     0.20601 water trees road river green long sun
27     0.04143 charles maria washington french great indians thou
28     0.02129 lake wylder rachel larkin stanley captain brandon
29     0.05588 sir wid good boy poor m'carthy rogue
30     0.59271 nan time good replied make made natter
31     0.07142 john mother big van flurry farn family
32     0.12677 replied god man father ha poor good
33     0.11696 sir marston mrs ma'am mademoiselle room de
34     0.22014 money man business good dollars thousand pay
35     0.03537 nartie sally billy lydia wallace manroe teddy
36     0.02801 sir toole puddock sturk dangerfield nutter doctor
37     0.01237 gerald ye mr ffrench young answered man
38     0.20356 mrs mother room girl woman girls margaret
39     0.02555 woodward harry replied sir barney mother goodwin
40     0.29075 nan head back dead hand nen horse
41     0.31523 time felt spirit eye heart eyes length
42     0.05723 school city paper editor wichta story job

Mar 28, 2017 8:56:35 AM cc.mallet.topics.ParallelTopicModel optimizeBeta
INFO: [beta: 0.01314]
Mar 28, 2017 8:56:35 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <400> LL/token: -9.24646
Mar 28, 2017 8:56:35 AM cc.mallet.topics.ParallelTopicModel estimate
INFO:
Total time: 1 minutes 6 seconds
.
```

Output screenshot - 10 topics instead of 43

```
Mar 28, 2017 9:01:07 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <350> LL/token: -9.20658
Mar 28, 2017 9:01:08 AM cc.mallet.topics.ParallelTopicModel optimizeBeta
INFO: [Data: 0.04899]
Mar 28, 2017 9:01:08 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <360> LL/token: -9.20642
Mar 28, 2017 9:01:09 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <370> LL/token: -9.20285
Mar 28, 2017 9:01:10 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <380> LL/token: -9.20313
Mar 28, 2017 9:01:12 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <390> LL/token: -9.20172
Mar 28, 2017 9:01:13 AM cc.mallet.topics.ParallelTopicModel estimate
INFO:
0      0.09582 irish men ireland washington french great people
1      0.23133 sir man mr replied good reilly country
2      0.16219 replied wld god man father good poor
3      0.28464 time mr man day years work good
4      0.13631 john mother big nell tin tin flurry farm
5      0.12586 ' sir 'i mr good lake poor
6      0.66484 night door room house back hand man
7      0.14838 susan marlie ers rachael miss mother billy
8      0.71019 heart father time love life day felt
9      0.05437 don d'arcy bryce man carigan bejober's parker

Mar 28, 2017 9:01:13 AM cc.mallet.topics.ParallelTopicModel optimizeBeta
INFO: [Data: 0.04045]
Mar 28, 2017 9:01:13 AM cc.mallet.topics.ParallelTopicModel estimate
INFO: <400> LL/token: -9.20225
Mar 28, 2017 9:01:13 AM cc.mallet.topics.ParallelTopicModel estimate
INFO:
Total time: 47 seconds
^I
```


More functionalities in Mallet

- ▶ `mallet.doc.topics` function returns topic weights for each document in the training data.
- ▶ `mallet.read.dir` function takes a folder path with `.txt` files and converts it into id-text format for `mallet`.
- ▶ `mallet.topic.hclust` function performs a clustering of topics.
- ▶ `mallet.topic.labels` function gives a string with most probable words for a topic.

More details:

<https://cran.r-project.org/web/packages/mallet/mallet.pdf>

Getting Topics and Words

doc.topics <- mallet.doc.topics(topic.model, smoothed=T, normalized=T)

```
> doc.topics[1,10]
[1] 0.000782777 0.000464610 0.000309616 0.002738675 0.003149773 0.000253809 0.173714346 0.000288194 0.003702725 0.704647811
> doc.topics
      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]      [,8]      [,9]     [10]
[1,] 0.238152988 0.000494647 0.000306316 0.000797400 0.008741511 0.000251138 0.134789945 0.002472086 0.00504650 0.335797-01
[2,] 0.017670471 0.000404567 0.000182885 0.024604682 0.155283660 0.002563680 0.158782983 0.000288986 0.004072931 0.034468-01
[3,] 0.000264295 0.000403856 0.000320653 0.157330280 0.119502470 0.002613796 0.237779794 0.000309228 0.06246802 0.188734-01
[4,] 0.011940862 0.000518276 0.000344121 0.000608078 0.007494326 0.000762512 0.343122524 0.000311435 0.000942371 0.948687-01
[5,] 0.0025183147 0.000518109 0.000332963 0.002870404 0.133511321 0.000773718 0.237224247 0.107670857 0.14633379 0.453090-01
[6,] 0.0002197372 0.0231217736 0.000362334 0.058585166 0.000651999 0.000294319 0.1446593257 0.0680033372 0.140422381 0.660273-01
[7,] 0.000295704 0.000500641 0.000332603 0.064062166 0.091549617 0.000759414 0.1049235792 0.0098115240 0.161468161 0.712694-01
[8,] 0.0002195791 0.000521638 0.000384006 0.074570625 0.000538300 0.000120061 0.064687244 0.011300576 0.00781132 0.315745-01
[9,] 0.0002138340 0.000520215 0.005231959 0.083661182 0.002736393 0.000266429 0.1163503995 0.0223832201 0.031018817 0.376862-01
[10,] 0.030738277 0.000466110 0.000306116 0.002738675 0.003149773 0.000253809 0.173714346 0.000288194 0.003702773 0.7046478-01
[11,] 0.018083141 0.000576237 0.0007842512 0.127306525 0.000487300 0.000515539 0.205632850 0.0711622375 0.11751398 0.374898-01
[12,] 0.000294314 0.012169592 0.000131263 0.147616140 0.000281224 0.005040791 0.225558082 0.0026424705 0.005650753 0.225141-01
[13,] 0.0002987361 0.0028518512 0.021794740 0.000692496 0.074182550 0.007431708 0.3090679774 0.0336917105 0.16377742 0.361336-01
[14,] 0.000297253 0.000584283 0.002703071 0.007700238 0.002652393 0.000777991 0.3707852014 0.000312688 0.059670781 0.161464-01
[15,] 0.0076717489 0.000529597 0.000350504 0.070079310 0.007753670 0.000291394 0.180497178 0.112019547 0.14560439 0.150763-01
[16,] 0.0002449489 0.0033914371 0.031178130 0.0518708062 0.050629650 0.000128078 0.2200082103 0.000307847 0.008675472 0.540542-01
[17,] 0.0002164774 0.025251262 0.000509487 0.124254763 0.094251328 0.000289953 0.243841337 0.0003268137 0.046220932 0.640945-01
[18,] 0.0002170140 0.005451210 0.000551816 0.00009179 0.000198493 0.0002980723 0.0361378037 0.1001363084 0.108479193 0.523494-01
[19,] 0.012120946 0.000587895 0.007408527 0.180485445 0.005853406 0.0002795847 0.101656340 0.0058832100 0.097024674 0.818250-01
[20,] 0.000942894 0.000520215 0.005231959 0.000624129 0.073562263 0.000266429 0.1163503995 0.000322829 0.031018817 0.621084-01
[21,] 0.0002159435 0.005438644 0.000338032 0.106831783 0.000297161 0.000755509 0.065704793 0.009653215 0.014082220 0.548292-01
[22,] 0.0001766907 0.0054075078 0.000327345 0.007623841 0.0504490185 0.000374692 0.268360725 0.0004245884 0.019383076 0.307186-01
[23,] 0.000214042 0.000532622 0.000589358 0.159956847 0.078495260 0.000296325 0.203917488 0.0053917074 0.001763933 0.488600-01
[24,] 0.000210665 0.025750415 0.000329648 0.048181200 0.107511519 0.0002693119 0.263555188 0.0003835479 0.008866796 0.662576-01
[25,] 0.0046093128 0.000245116 0.000338518 0.031194627 0.003764659 0.000773074 0.2017145107 0.002363174 0.00008837 0.752626-01
[26,] 0.000214042 0.000532622 0.000589358 0.048181200 0.107511519 0.0002693119 0.263555188 0.0003835479 0.008866796 0.662576-01
[27,] 0.000218951 0.0003021526 0.000353571 0.140128901 0.017360798 0.027180314 0.0713412127 0.002567324 0.191048596 0.331842-01
[28,] 0.000193337 0.000458217 0.02587677 0.161458406 0.000729403 0.000259516 0.101120740 0.041878305 0.10571835 0.680802-01
[29,] 0.000218951 0.0003021526 0.000353571 0.140128901 0.017360798 0.027180314 0.0713412127 0.002567324 0.191048596 0.331842-01
[30,] 0.000214375 0.027452187 0.000475119 0.007970940 0.002743293 0.000287141 0.2072289419 0.001125118 0.104843757 0.118904-01
[31,] 0.000217749 0.000315519 0.000303027 0.000181470 0.00013533 0.009468140 0.000181351 0.011830821 0.000181351 0.790474-01
[32,] 0.000214796 0.000520620 0.000616474 0.075091040 0.0003815794 0.000304461 0.0051308792 0.010720802 0.001479139 0.742118-01
[33,] 0.000159851 0.0002908814 0.0003138256 0.004974120 0.057683170 0.002463903 0.000291408 0.025766409 0.063858-01
[34,] 0.000197854 0.007471252 0.01381405 0.000577593 0.027936813 0.000650951 0.4444299330 0.0002387455 0.033218141 0.709705-01
```

Sensitivity of a topic model to its parameters

- ▶ Note: Topic models are very sensitive to their parameters.
- ▶ How to choose the number of topics?: Intuitive - needs an understanding of the corpus. There are some ways of using probability based estimates (which are beyond the scope of this course), but there is no "formula" for choosing the best number.

Sensitivity of a topic model to its parameters

- ▶ Note: Topic models are very sensitive to their parameters.
- ▶ How to choose the number of topics?: Intuitive - needs an understanding of the corpus. There are some ways of using probability based estimates (which are beyond the scope of this course), but there is no "formula" for choosing the best number.
- ▶ How to choose the number of iterations?: Looking at log likelihoods is a good way.

Useful reference: ldatuning library in R

<https://cran.r-project.org/web/packages/ldatuning/ldatuning.pdf>

Human evaluation of topic models -1

- ▶ Word intrusion test: If an irrelevant word is shown along with related words for a given topic, a human evaluator should be able to identify that word as an intruder.

Human evaluation of topic models -1

- ▶ Word intrusion test: If an irrelevant word is shown along with related words for a given topic, a human evaluator should be able to identify that word as an intruder.
- ▶ If a human is having a difficulty doing that, the topic is not coherent.

Human evaluation of topic models -1

- ▶ Word intrusion test: If an irrelevant word is shown along with related words for a given topic, a human evaluator should be able to identify that word as an intruder.
- ▶ If a human is having a difficulty doing that, the topic is not coherent.
- ▶ Experimental results: models with a larger number of topics may improve automated evaluations but will reduce human interpretability, because they will not be coherent.

Human evaluation of topic models -2

- ▶ Topic intrusion test: If a document is given, and its most appropriate topics (according a topic model) are shown along with a random topic, human evaluator should be able to identify the intruder topic correctly.

Human evaluation of topic models -2

- ▶ Topic intrusion test: If a document is given, and its most appropriate topics (according a topic model) are shown along with a random topic, human evaluator should be able to identify the intruder topic correctly.
- ▶ If topic assignments from the model are intuitive, then the human intruder should not have difficulties in doing this.

Human evaluation of topic models -2

- ▶ Topic intrusion test: If a document is given, and its most appropriate topics (according a topic model) are shown along with a random topic, human evaluator should be able to identify the intruder topic correctly.
- ▶ If topic assignments from the model are intuitive, then the human intruder should not have difficulties in doing this.
- ▶ topic here means top-N keywords from the topic model for that topic.

Human evaluation of topic models -2

- ▶ Topic intrusion test: If a document is given, and its most appropriate topics (according a topic model) are shown along with a random topic, human evaluator should be able to identify the intruder topic correctly.
- ▶ If topic assignments from the model are intuitive, then the human intruder should not have difficulties in doing this.
- ▶ topic here means top-N keywords from the topic model for that topic.
- ▶ Experimental results: humans do well documents which have focussed discussions since it is easy to assign a topic in such cases (easy for machines too!)

Automatic evaluation of topic models-1

1. Approximating the word intrusion test automatically by comparing semantic relatedness between words in a topic.

Automatic evaluation of topic models-1

1. Approximating the word intrusion test automatically by comparing semantic relatedness between words in a topic.
2. hold out a few words from each document from training data and use it to evaluate the topic model.

Note: Inference from a topic model is available in original Mallet tool, but does not seem to be available in the R library version.

<http://mallet.cs.umass.edu/topics.php>

Enthusiasts can have a look at this R code for the same:

<https://gist.github.com/agoldst/edcfd45b5ac371296b76>

Automatic evaluation of topic models -2

- ▶ Use the topic model as a part of some other application scenario (e.g., using it to retrieve similar documents to a given document).
- ▶ Check if the application performance has gotten better because of the use of this topic model.
- ▶ If it does, we say the topic model is good for that application. Else, not good.

Automatic evaluation of topic models -2

- ▶ Use the topic model as a part of some other application scenario (e.g., using it to retrieve similar documents to a given document).
- ▶ Check if the application performance has gotten better because of the use of this topic model.
- ▶ If it does, we say the topic model is good for that application. Else, not good.
- ▶ Semi-automatic: comparing topics in terms of overlapping words in Top-25 words for each topic.

Evaluation: general remarks

- ▶ A general conclusion from topic modeling research has been that human evaluations and machine evaluations do not agree with each other.
- ▶ While machine evaluations are good for certain tasks (e.g., search engines etc.) we may need topic models that are optimized to human interpretations for some other tasks (e.g., analyzing literary documents for themes).

Evaluation: general remarks

- ▶ A general conclusion from topic modeling research has been that human evaluations and machine evaluations do not agree with each other.
- ▶ While machine evaluations are good for certain tasks (e.g., search engines etc.) we may need topic models that are optimized to human interpretations for some other tasks (e.g., analyzing literary documents for themes).
- ▶ One thing to keep in mind: human evaluation is expensive.
- ▶ How to consider human interpretation as a part of the mathematical modeling process? - is an ongoing research question.

Next Class

- ▶ Wrapping up topic modeling
- ▶ Topic modeling with or withoutallet - Assignment 6 practice