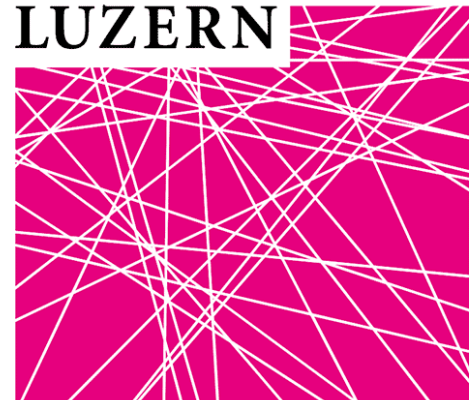


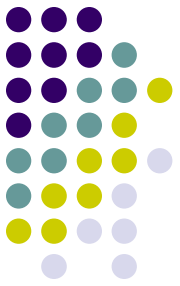
Big Data Analytics

Lecture 4/A Semisupervised classification methods



UNIVERSITÄT
LUZERN

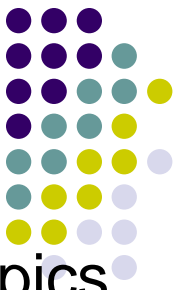




Reference

- ✓ Kohei Watanabe and Yuan Zhou (2020) Theory-Driven Analysis of Large Corpora: Semisupervised Topic Classification of the UN Speeches, *Social Science Computer Review*, DOI: 10.1177/0894439320907027
- ✓ Shusei Eshima, Kosuke Imai, Tomoya Sasaki (2020). Keyword Assisted Topic Models, arXiv:2004.05964v1

Semisupervised classification



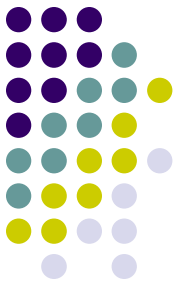
When using **topic models**, most researchers have the topics of interest in mind and possess (hopefully) a substantial amount of knowledge about them

After all, social scientists analyze textual data in order to empirically test hypotheses derived from substantive theories

Thus, researchers should find it natural to incorporate such **prior information** into topic models as keywords

In contrast, the standard topic models discussed so far are designed for the settings in which researchers wish to **explore** the contents of corpus, w/o any prior knowledge (with the exception of some priors on the covariates that could affect for example the distribution of thetas in a STM)

Semisupervised classification



Using *unsupervised classification algorithms* produces therefore results that sometimes are difficult to make sense of (and perhaps not the topics you would be interested in as a researcher...)

And so?

Perhaps we could give a try to **semisupervised classification**

Semisupervised classification



How do they work? Let's first discuss **Newsmap** (but the logic, as we will see, generalizes also to other algorithms within this approach)

A four-steps procedure

First step: you define a set of pre-defined categories in which you want to classify texts

For example, let's suppose you want to classify the news in our corpus according to two class-labels: either if they discuss about Ukraine or Iraq

Semisupervised classification



Second step: you have to identify some words that would help us later on to *automatically* create a dictionary

Let's call these special words «**seed words**»

In our case, the seed words for countries such as Ukraine and Iraq could be {Ukraine, Ukrainian*, Kiev} and {Iraq, Iraqi*, Baghdad}

Note: the list of seed words is the only manual input to the system (beyond the definition of the list of categories you are interested about), and serves as **semi-supervision**

Semisupervised classification



Third step: the training-stage

In the **training-stage** the algorithm takes advantage of these seed-words to **automatically create a dictionary** for each label. How to do that?

Firstly, it searches individual documents for keywords in the seed dictionary (simple keyword matching) and gives them class labels (countries)

For example, suppose that in the corpus we have the following text: “This is an article about *Ukraine*”

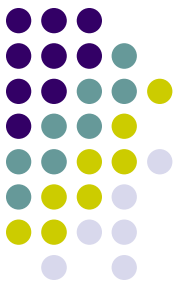
The algorithm will **automatically label** such document as “Ukraine” given that it includes at least one of the seed-word {Ukraine, Ukrainian*, Kiev}

Semisupervised classification

Secondly, we aggregate the frequency of words according to the class labels to create **contingency tables**, i.e., the labels are used to estimate the association between the labels and features



Semisupervised classification



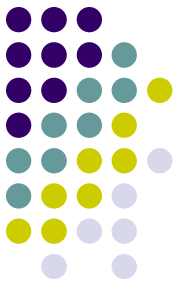
In the contingency table below, c_j is a country of interest (i.e., Ukraine) and \bar{c}_j refers to all the other countries (i.e., Iraq in our example)

w_i is the word for which scores are calculated (say word “article”) and w'_i refers to all the other words

F values are all raw frequency counts of words in respective classes for different subset of texts (i.e., c_j : all those texts that report at least one of the seed-words we defined with respect to Ukraine; \bar{c}_j : all those texts that report at least one of the seed-words we defined with respect to Iraq)

	c_j	\bar{c}_j
w_i	F_{11}	F_{01}
w'_i	F_{10}	F_{00}
$w_i + w'_i$	$F_{1.}$	$F_{0.}$

Semisupervised classification



The estimated score \hat{S} of word w_i for class-label (country in our case) c_j is then estimated by focusing on **co-occurrences of words**

In particular \hat{S}_{ij} (i.e., the “association score of word i for label j ”) is calculated as the association between w_i and c_j subtracted by the association between w_i and \bar{C}_j :

$$\hat{S}_{ij} = \log \frac{F_{111}}{F_{1.}} - \log \frac{F_{01}}{F_{0.}}$$

Semisupervised classification



$$\hat{S}_{ij} = \log \frac{F_{11j}}{F_{1.}} - \log \frac{F_{01j}}{F_{0.}}$$

For example, if the word “`article`” appears 4 times in the subset of texts labelled as c_j and this subset includes 100 words, while it appears only once in the subset of texts labelled as \bar{c}_j and this subset includes 300 words, then...

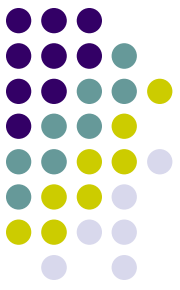
the score for “`article`” with respect to c_j will be: $\log 4/100 - \log 1/300 = +1.08$

and with respect to \bar{c}_j it will be: $\log 1/300 - \log 4/100 = -1.08$

As a result, a text that includes only the word “`article`” will be classified, *ceteris paribus*, as belonging to c_j

Of course, you have to replicate this exercise to estimate a score for all the words in those texts that include at least one seed-word

Semisupervised classification



This table shows fictional scores given to the words most strongly associated with Ukraine and Iraq

	Ukraine	Score	Iraq	Score
1	Ukraine	11.84	Iraq	11.58
2	Ukrainian	10.36	Baghdad	10.56
3	Kiev	10.34	Iraqi	10.39
4	Ukrainians	7.94	Iraqis	8.15
5	Poroshenko	7.64	Anbar	8.14
6	Mariupol	7.15	Ramadi	7.55
7	Yatseniuk	6.94	Fallujah	7.51
8	Donetsk	6.92	Falluja	7.50
9	Slovyansk	6.84	Kirkuk	7.42
10	Lugansk	6.72	Tikrit	7.36

Through this approach, many new words are identified based on **co-occurrences** both for Ukraine and Iraq! These words will be added to the **dictionary** you created (with a given estimated weight) **starting** from the original list of seed-words

Semisupervised classification



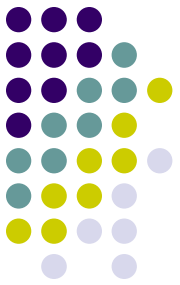
Fourth step: The classification-stage

Newsmap then predicts the class-label (i.e., \hat{C} - countries in our case) most strongly associated with documents simply by finding a class-label that yields the largest total scores \hat{S} weighted by the normalized frequency of word f_i in documents:

$$\hat{C} = \arg \max_j \sum_i \sum_j \hat{s}_{ij} f_i.$$

i.e., a given document j will be assigned to Ukraine rather than Iraq if it presents words most frequently used in the subset of documents that use one of the seed-words related to Ukraine {Ukraine, Ukrainian*, Kiev} , even if it does not include **any of these seed-words!**

Semisupervised classification: the 4 steps



1. The researcher pre-defines a list of labels/categories in which she is interested (Ukraine and Iraq in our example)



2. The researcher identifies the **seed words** associated to each pre-defined categories (the only manual input to the system)



3. In the **training-stage** the algorithm takes advantage of these seed-words to automatically create a dictionary for each label



4. Then, in the **classification-stage** the algorithm classifies all the texts (also texts that do not include *any* of our seed-words) into one of the pre-defined categories

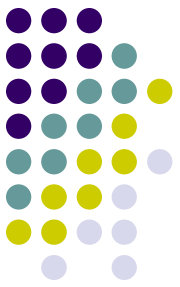
Semisupervised classification



This approach is advantageous because:

- (1) contrary to a fully unsupervised model, you **know ex-ante** the label of the topics you are looking for!
- (2) training of new models **does not require** any manual coding but only a **seed word dictionary** used to train semisupervised document classifiers
- (3) creating a seed word dictionary is relatively easy because the number of words required for a seed word dictionary is a **fraction of a standard dictionary approach**
- (4) dictionaries trained via seed-words can then be ported to **different projects** without or little modification

Semisupervised classification



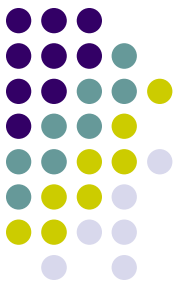
However be careful...

- ✓ *First*: use of semisupervised models requires a good knowledge of what you are analyzing!
- ✓ This is going to affect the specific selection of the categories you want to investigate given your corpus and (*crucially*) the related seed-words

In particular, seed-words should be of **high quality**!

1. Seed words should be *knowledge-based*, i.e., words selected based on researchers' background knowledge in the field
2. Seed-words should be (reasonably) present in your corpus (otherwise the training-stage will perform poorly...)

Semisupervised classification

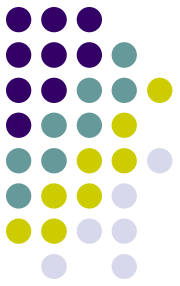


Watanabe and Zhou (2020) also suggests to eventually choose seed-words from the most frequent words in the corpus (i.e., *Frequency-based seed words*)

While acknowledging that *knowledge-based seed words* are superior to *frequency-based seed words* because knowledge-based seed words offer operational definitions of the concepts and have greater external validity that ensures portability across corpora...

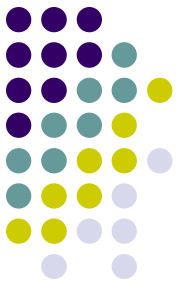
...researchers can still collect a set of words that are frequent in the corpus to increase the seed-word coverage that can be used to train the semisupervised machine learning model

Semisupervised classification



Eshima et al. (2020) argue *however* to be cautious about this latter option, because that would imply analyzing the same data as the one used to derive the seed words (i.e., endogeneity problems)

Semisupervised classification

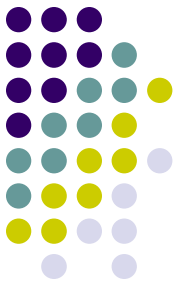


For several suggestions on how to improve a seed-words dictionary list see: Watanabe and Zhou (2020)

For example, they suggest as a rule-of-thumb that the optimal size of seed dictionary is approximately six to seven words for each topic

But of course, this also depends on the specific analysis to be considered

Semisupervised classification



- ✓ *Second:* Newsmap also requires users to define **all the relevant topics** in a seed dictionary because it estimates features' association with a topic by comparing between their frequencies in documents with the label and all other labels, **ignoring documents without labels**

Semisupervised classification



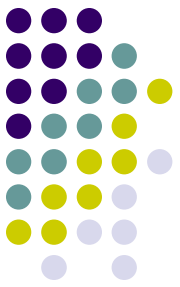
So, for example, if in your corpus you have news discussing about Ukraine, Iraq, but also Italy, then if you have not identified seed-words about Italy as well, Newsmap will ignore such documents in its training-stage!

As a result, in the classification stage, the “true” texts discussing about Italy will be either:

- a) classified as Iraq;
- b) classified as Ukraine;
- c) classified as NAs if they do not include any of the words that are present in the subset of texts that include at least one of your seed-word

Under a) and b) the possibility of false positives clearly can explode!

Semisupervised classification



An alternative is employing a **semi-supervised topic model** (Eshima et al. 2020; for an application: Curini and Vignoli 2020)

The main idea: we assume that in our corpus there are two types of topics

✓ **topics with keywords** (or seed-words using the jargon of Newsmap) defined ex-ante by the researcher, which are of primary interest to researchers and are referred to as *keyword topics*

For example, going back to our previous discussion, we can have 2 *keyword topics*: a topic related to Ukraine defined by the keywords {Ukraine, Ukrainian, Kiev} and a topic about Iraq defined by the keywords {Iraq, Iraqi, Baghdad}

Semisupervised classification



And..

✓ **topics without keywords**, called *no-keyword topics*

One key-difference with Newsmap is that you can have therefore one or more of off-keywords topics (beyond the ones you identify via your seed-words selection)

And you can still recover them from the analysis: this would by definition decrease the risk of false-positives!

Plus: this approach also assumes that each text is a *mixture of topics*, i.e., if we have a text discussing about Ukraine *and* Iraq as well, we are able to recover it from the analysis

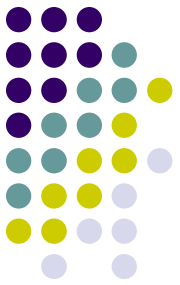
Semisupervised classification



How it works? Going back to our example, let's suppose that we select $k=3$, with 2 keyword topics (Ukraine and Iraq) and 1 no-keyword topic

A semi-supervised topic model is a traditional topic model with a crucial difference in the **third step of the procedure** (i.e., after having extracted θ_i and β_k)

Classification methods

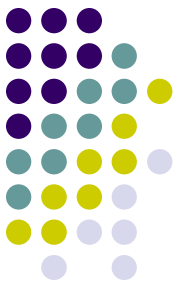


1. Choose $\theta_i \sim \text{DIRICHLET}(\alpha)$
2. Choose $\beta_k \sim \text{DIRICHLET}(\delta)$
3. Choose a topic $z \sim \text{Multinomial}(\theta_i)$

If the sampled topic is one of the no-keyword topics, then we follow our usual approach. That is, we draw word w_i from the corresponding word distribution of the topic.:

- Choose a word $w_i \sim \text{Multinomial}(\beta_{i,k=z})$

Classification methods

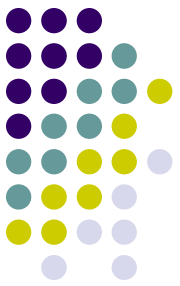


1. Choose $\theta_i \sim \text{DIRICHLET}(\alpha)$
2. Choose $\beta_k \sim \text{DIRICHLET}(\delta)$
3. Choose a topic $z \sim \text{Multinomial}(\theta_i)$

If *however* the sampled topic is a keyword topics (i.e., either Ukraine or Iraq)...

- ✓ we first draw a Bernoulli random variable with success probability p for word w in document i
- ✓ If the Bernoulli random variable is equal to 0, then we sample the word from the standard topic-word distribution of the topic
- Choose a word $w_i \sim \text{Multinomial}(\beta_{i,k=z})$
- ✓ **In contrast**, if this variable is equal to 1, then word w is drawn from the set of keywords for the topic that we defined ex-ante

Semisupervised classification

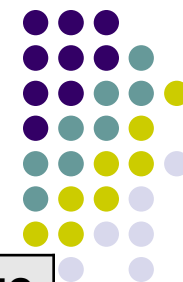


Let's see an example. Suppose you have these 3 sentences in your corpus with the following DfM:

	world	travel	dangerous	statistics	Kiev	Baghdad	luggage
Doc1	4	1	0	2	1	0	2
Doc2	0	2	3	1	2	2	2
Doc3	2	1	2	2	0	0	3

Now you decide to fit a semi-supervised topic model with the previous 2 keyword-topics – K1 and K2 (with the list of keywords: {Ukraine, Ukrainian, Kiev} and {Iraq, Iraqi, Baghdad}), and 1 no-keyword topic – NK1

Classification methods



Step 1 to Step 3 – first random assignment

	world	travel	dangerous	statistics	Kiev	Baghdad	luggage
Doc1	4	1	0	2	1	0	2
Doc2	0	2	3	1	2	2	2
Doc3	2	1	2	2	0	0	3



	K1	K2	NK1
Doc1	0.6	0.0	0.4
Doc2	0.2	0.4	0.4
Doc3	0.4	0	0.6



	world	travel	dangerous	statistics	Kiev	Baghdad	luggage
K1	0.1	0.2	0.1	0.1	0.3	0.0	0.2
K2	0	0.2	0.2	0.1	0.1	0.3	0.1
NK1	0.1	0.2	0	0.4	0	0	0.3

Document-topics matrix
(first assignment) - Step 1



Topic-terms matrix (first
assignment) – Step 2

Step 3 – suppose you draw NK1 for document *Doc1*. You know from Step 1 that 40% of words should be devoted to NK1 – i.e., 4 out of 10 words from *Doc1*, and 60% to K1. We can now generate the document *Doc1* using our usual approach

Classification methods



That is...when generating the document Doc1 (consisting of 10 words):

When we sample NK1, we follow the typical topic-model approach

We pick the first word (out of 4) to come from the beta distribution of NK1 topic, which then gives you the word “statistics”

We pick the second word to come from NK1, which gives you “world”

We pick the third word to come from NK1, giving you “luggage”

We pick the fourth word to come from NK1, giving you “travel”

Classification methods



Now you sample $K1$ (a keyword topic!). 6 words should be related to it.

For the first word (out of 6) you draw a Bernoulli distribution = 0. As a result you pick the word to come from the beta distribution of $K1$ topic, which gives you the word “dangerous”

For the second word, Bernoulli = 1. As a result you randomly pick the word to come *from the keywords of $K1$ topic*, which gives you the word “Ukraine”

For the third word, Bernoulli = 1. Once again, you pick the word to come *from the keywords of $K1$ topic*, which gives you the word “Ukrainian”

For the fourth word, Bernoulli = 0. As a result you pick the word to come from the beta distribution of $K1$, which gives you the word “travel”

For the fifth word, Bernoulli = 0. As a result you pick the word to come from the beta distribution of $K1$, which gives you the word “world”

For the sixth word, Bernoulli = 1. As a result you pick the word to come *from the keywords of $K1$ topic*, which gives you the word “Ukraine”

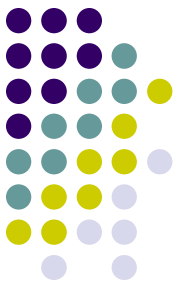


Classification methods

So the document Doc1 generated under our assumption will be “*statistics world luggage travel dangerous Ukraine Ukrainian travel world Ukraine*” (in a bag-of-words approach)

Then we can do the same for Doc2 and Doc3

Starting from this first assignment, we try to improve our model following substantially the same strategy discussed with respect to a Topic model

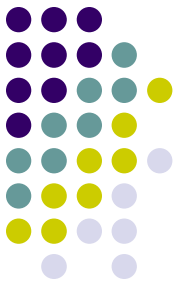


Classification methods

The semi-supervised topic is therefore based on a *mixture of two distributions*, one with positive probabilities only for keywords and the other with positive probabilities for all words

It is straightforward to show that this mixture structure makes the prior means for the frequency of user-selected keywords given a topic greater than those of non-keywords in the same topic (as can be seen above)

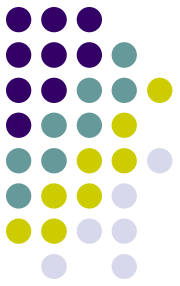
Semisupervised classification



In addition, the prior variance is also larger for the frequency of keywords given a topic than for non-keywords

This encourages the algorithm to place *a greater importance on keywords a priori* while allowing the model *to learn from the data* about the precise degree to which keywords matter for a given topic

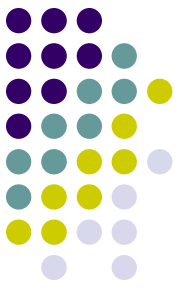
Semisupervised classification



After having fitted a semisupervised topic model, each topic with keywords already has a label and so there is no need to interpret the resulting topics after model fitting (one topic is going to be Ukraine and the other is going to be Iraq) – of course this happens as far the keywords you have selected are of good quality (more on this below)

In contrast, the no-keywords topics require as usual a post-hoc labeling (could it be related to Italy?)

Semisupervised classification



Note that the algorithm we will employ is generally robust to the number of no-keywords topics (as long as they are >0)

Why that? Cause if the keywords are of high quality (more on this below...) and they identify as 20% of document i as discussing about keyword topic X , it does not matter if the remaining 80% of document i is related to 3 or 10 no-keywords topic; the 20% devoted to topic X will remain (approximately) still there

This of course also implies that everytime you are estimating a semi-supervised topic model, you are mainly interested in retrieving information about the keyword-topics, while the no-keyword topics are usually treated as simple “noise”

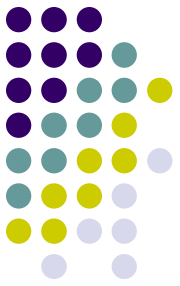
Semisupervised classification



The selection of **high quality keywords** is once again critical for the successful application of such approach!

1. The keywords selection should be theoretically sound!
That also implies that (within each topic) keywords should refer to the same topic!
2. The keywords should be present in a non-negligible way in the corpus (this is something you can **check ex-ante**)
3. Moreover, the unique keywords selected for each topic should appear frequently in the keyword topic, i.e., they should discriminate their topics from others (this is something you can **check only ex-post**)

Semisupervised classification



Interestingly, you can incorporate within a semi-supervised topic model *also covariates* for the document-topic distribution as social scientists often have meta information about documents

This is similar to what we discussed for STM (so we are not going to enter into that!). Do you remember?

We can have, in other words, a semisupervised *structural* topic algorithm!