

# Mining for spillovers in patents

## Plan of attack

Jacob Calder<sup>\*</sup>, Iwo Dubaniowski<sup>†</sup>, C. Vic Hu<sup>‡</sup>, Subhashini Venugopalan<sup>§</sup>

EE380L Data Mining  
University of Texas at Austin  
<https://github.com/cvhu/ee380l-ghosh-project>

March 25, 2013

## 1 Approach

Here is a detailed systematic approach to make inferences from the topic model.

1. Consolidate the data consisting of - abstract, claims, date of filing/grant, geographical location, technology area, classification numbers, patent-holder information and other (less important) details.
2. Run the topic model (a few times) on the abstract and claims, to generate a list of consistent topics word clusters. Call this *TopicWordsCluster<sub>TM</sub>*.
3. To make any inference, we first need to associate the *TopicWordsCluster<sub>TM</sub>* with the actual topic (e.g. with respect to the example in Topic Modelling web page and tutorial (<http://www.programminghistorian.org/lessons/topic-modelling-and-mallet>), that would be cricket, movies, wildlife etc.) And we don't know these topics, so we need to identify them based on the technology areas and classification numbers.

- (a) Identify the technology areas correlated with each *TopicWordsCluster<sub>TM</sub>*.

This can be done by first identifying the most correlated topics for each document. Then, associating the document's technology area with that *TopicWordsCluster<sub>TM</sub>*. Doing this for all documents, we can get a list of technology areas corresponding to each topic, and counting the frequencies will give us a single technology area for each topic.

This has been tested on the abstracts. The scripts and results can be found in <https://github.com/cvhu/ee380l-ghosh-project/tree/master/data/pyAnalysisOutput>

---

<sup>\*</sup>jcalder@utexas.edu

<sup>†</sup>mduban@gmail.com

<sup>‡</sup>vic@cvhu.org

<sup>§</sup>vsubhashini@utexas.edu

- (b) Similarly associate classification classes with every *TopicWordsCluster<sub>TM</sub>*.  
 To accomplish this, first get the classification names and details corresponding to each of the classification numbers for each patent. Now, the algorithm is very similar to the previous step (as in technology areas). For each patent, identify the most correlated topic and associate it's classification details to that topic. Then refine by topic to identify the list of classifications for each *TopicWordsCluster<sub>TM</sub>*, and pick the top 3-4 classifications.
4. From the previous step, we should have now generated a set of *Topics<sub>TM</sub>* or topic names for each *TopicWordsCluster<sub>TM</sub>* based on technology area and patent classification details. The next step is to do the actual inference with respect to geographies and time.
- (a) For geographies (and time), run the topic model again to generate a bunch of *TopicWordsCluster<sub>Geo</sub>* (correspondingly *TopicWordsCluster<sub>Years</sub>*). Map *TopicWordsCluster<sub>Geo</sub>* (*TopicWordsCluster<sub>Years</sub>*) with *TopicWordsCluster<sub>TM</sub>* based on the words directly (or re-doing step 3). The result of this will be *Topics<sub>Geo</sub>* (and *Topics<sub>Years</sub>*).
- (b) Now we will have for each geography the set of most frequent *Topics*, and we should be able to test our hypothesis:
- Are certain geographies (and years) associated with specific topics?
  - Is there a progression in the set of topics across years?
  - Map *Topics<sub>Geo</sub>* (and *Topics<sub>Years</sub>*) to see if (*Topic<sub>i</sub>, Geo<sub>j</sub>, Year<sub>k</sub>*) has some association with (*Topic<sub>i</sub>, Geo<sub>l</sub>, Year<sub>j+1</sub>*) etc.