The Hard Problem of Prediction for Prevention Reading Between the Lines

Hannes Mueller IAE (CSIC) Christopher Rauh (University of Montreal

Fundació Economia Analítica

April 2018

Financial Support from Fundación BBVA

- Civil wars are a serious humanitarian and economic problem.
- And we fail to prevent them.
- This is reflected in large expenditures on crisis response.
 - Humanitarian response: ca. 24.5 billion US dollars in 2014
 - Peacekeeping: ca. 8 billion US dollars per year.
- Review of the United Nations Peacebuilding Architecture (2015):
 - If more global priority were consistently given to efforts at sustaining peace, might there not, over the course of time, be reduced need for crisis response?

- Define hard problem:
 - why hard
 - why problem
- Literature
- Solution:
 - $\bullet\,$ topic model with news text $\rightarrow\,$ summaries of text
 - use summaries to get early warning
 - speculate why it works

- Violence data from Uppsala Conflict Data Program (UCDP)
- Georeferenced Event Dataset (GED), Sundberg and Melander (2013)
- Gives quarterly data for all countries 1989-2016
- Internal conflicts: state-based conflict, non-state conflict, one-sided violence.
- We focus on onset of conflict, i.e. code a dummy of start.
- Conflict in the literature is defined as 25+ or 1000+ battle-related deaths per year.
- Not obvious how to translate this to the quarterly data.
- We use three thresholds: 1, 50 and 500 (all violence, top 50% and top 25%)

- Take the 50+ definition.
- There were 433 onsets in almost 19,000 observations.
- Conflict history is a strong predictor of conflict onset.
- 359 onsets followed within 10 years of another conflict, 75 were "new" onsets.
- The following plot shows the risk of an onset post-conflict.

The Hard Problem of Prediction

Onset likelihood in post-conflict period (50+):



- Average risk outside the 10-year period is very low: about 0.6 percent.
- This means post-conflict period is a powerful forecast.
- After around 40 quarters, however, risk is again close to "normal".
- Call the 40 quarters after conflict post-conflict.
- Call all other quarters with peace pre-conflict.

• Then we can write down a simple Markov transition matrix:

		this quarter		
		peace pre- conflict	conflict	peace post- conflict
next quarter	peace pre- conflict	99.45%	0.00%	1.45%
	conflict	0.55%	80.21%	8.26%
	peace post- conflict	0.00%	19.79%	90.29%

• Peace always follows peace but is much less stable post-conflict .:

		this quarter		
		peace pre- conflict	conflict	peace post- conflict
next quarter	peace pre- conflict	99.45%	0.00%	1.45%
	conflict	0.55%	80.21%	8.26%
	peace post- conflict	0.00%	19.79%	90.29%

- Still, we have 75 onsets pre-conflict.
- This is what we call the hard problem.

The Hard Problem: Why Should We Care?

- Ironically, because it is hard.
- Pre-conflict peace is stable.
- Post-conflict peace is unstable.
- After 30 quarters, the distribution when starting from the pre-war peace is

```
(0.86, 0.05, 0.09)
```

- We have an 86 percent likelihood to be in pre-war peace and a 5 percent chance to be in conflict.
- But when starting from post-war peace the distribution is

(0.26, 0.23, 0.51)

• We have a 51 percent likelhood to be in post-war peace and a 23 percent chance to be in conflict.

- High risk post-conflict means that it is a bad state to be in.
- If you prevent an escalation into conflict pre-conflict you prevent a country from entering a bad cycle.
- This makes prediction pre-conflict particularly important.

- Define hard problem:
 - why hard
 - why problem
- Literature
- Solution:
 - $\bullet\,$ topic model with news text $\rightarrow\,$ summaries of text
 - use summaries to get early warning
 - speculate why it works

Literature (Practise)

Standard fragility measure from the fund for peace.



Drivers of conflict:

- Ethnicity (Montalvo and Reynal-Querol (2005, AER), Esteban et al (2014, AER), Michalopoulos and Papaioannou (2016, AER),
- Weather (Miguel et al (2004, JPE), Hsiang et al (2013, Science), Ciccone (2011, AEJApplied), Searson (2015, JDE))
- Commodities (Bazzi and Blattman (2014, AEJMacro), Berman et al (2017, AER))
- Political Institutions (Besley and Persson (2011, QJE))

• Forecasts:

- Goldstone et al (2010, AJPS), Chadefaux (2012, JPR),
- ICEWS, Ward et al (2011)
- Impossibility of perfect forecast: Chadefaux (2017, JCR), Cedermann and Weidmann (2017, Science)

- Mueller and Rauh (forthcoming, APSR)
- Country fixed effects as a solution for the hard problem.
- For forecast, use only within variation.
- Advantage: you can be sure you predict the timing.
- Disadvantage: you throw away meaningful variation.
- Will try to solve the hard problem differently.

- Mullainathan and Spiess (2017, JEP)
- Machine learning is great in forecasting: put in x and get out \hat{y}
- In economics we are generally interested in hypothesis testing.
- But: useful for heterogenous data (text, images, recording, video).
 - Donaldson and Storeygard on images (2016, JEP)
 - Gentzkow et al on text (2017, NBER)
- We will use it in two ways:
 - unsupervised (feature extraction)
 - supervised (forecast)

- 700,000 articles from New York Times, Washington Post and Economist
- 3.7 million articles from BBC Monitor
- BBC Monitor tracks broadcast, press and social media sources in multiple languages from over 150 countries worldwide.
- Journalists filter, translate and report breaking news.
- We download an article if a country name or capital name is in the title.
- about 4.4 million articles dated from 1989q1 to 2017q3 on over 190 countries.

Number of Articles over Time



- We treat each article m as a vector of tokens w_m (police, bank, american president, united nations...)
- After dropping rare tokens we have 0.8 million tokens. Need to reduce dimensionality.
- Latent Dirichlet allocation (LDA) introduced by Blei, Ng, and Jordan (2003).
- Topics: probability distributions over the tokens.
- Text generation: journalist picks topic randomly then randomly picks tokens.
- "Latent": only the \mathbf{w}_m are actually observed.
- We tried K = 5, 10, 15 topics: low number of topics.
- The following pictures visualize our procedure.

Example: NYT - March 29, 1991. Libya

The exiled Prince Idris of Libya has said he will take control of a dissident Libyan paramilitary force that was originally trained by American intelligence advisers, and he has promised to order it into combat against Col. Muammar el-Qaddafi, the Libyan leader. The United States' two-year effort to destabilize Colonel Qaddafi ended in failure in December, when a Libyan-supplied guerrilla force came to power in Chad, where the original 600 commandos were based. The new Chad Government asked the United States to fly the Libyan dissidents out of the country, beginning a journey that has taken them to Nigeria, Zaire and finally Kenya. So far, no country has agreed to take them permanently. The 400 remaining commandos, who have been disarmed, were originally members of the Libyan Army captured by Chad in border fighting in 1988. They volunteered for the force as a way of escaping P.O.W. camps. "Having received pledges of allegiance from leaders of the force. Prince Idris has stepped in to assume responsibility for the troops' welfare," said a statement released in Rome by the royalist Libyan government in exile. It was overthrown in 1969.

Example: NYT - March 29, 1991. Libya (Stopwords)

the exiled prince idris of libya has said he will take control of a dissident libyan paramilitary force that was originally trained by american intelligence advisers, and he has promised to order it into combat against col. muammar el-qaddafi, the libyan leader. the united states' two-year effort to destabilize colonel gaddafi ended in failure in december, when a libyan-supplied guerrilla force came to power in chad, where the original 600 commandos were based. the new chad government asked the united states to fly the libyan dissidents out of the country, beginning a journey that has taken them to nigeria, zaire and finally kenya. so far, no country has agreed to take them permanently. the 400 remaining commandos, who have been disarmed, were originally members of the libyan army captured by chad in border fighting in 1988. they volunteered for the force as a way of escaping p.o.w. camps. "having received pledges of allegiance from leaders of the force, prince idris has stepped in to assume responsibility for the troops' welfare," said a statement released in rome by the royalist libyan government in exile. it was overthrown in 1969.

Example: NYT - March 29, 1991. Libya

exiled prince idris libya control dissident libyan paramilitary force originally trained american intelligence advisers, promised order combat col. muammar el-qaddafi, libyan leader. united states' two-year effort destabilize colonel gaddafi ended failure december, libyan-supplied guerrilla force came power chad, original 600 commandos based. new chad government asked united states fly libyan dissidents country, beginning journey taken nigeria, zaire finally kenya. far, country agreed permanently. 400 remaining commandos, disarmed, originally members libyan army captured chad border fighting 1988. volunteered force wav escaping p.o.w. camps. "having received pledges allegiance leaders force, prince idris stepped assume responsibility troops' welfare," statement released rome royalist libyan government exile. overthrown 1969.

Example: NYT - March 29, 1991. Libya (Lemmatizing)

exiled prince idris libya control dissident libyan paramilitary force originally trained american intelligence advisers, promised order combat col. muammar el-qaddafi, libyan leader. unit state two-year effort destabilize colonel gaddafi ended failure december, libyan-supplied guerrilla forces came power chad, origin 600 commando based. new chad government asked united states fly libyan dissidents country, beginning journey taken nigeria, zaire finally kenya. far, country agreed permanently. 400 remain commandos, disarmed, originally members libyan army captured chad border fighting 1988. volunteered force way escaping p.o.w. camps. "having received pledges allegiance leader force, prince idris steped assume responsibility troop's welfare," statement released rome royalist libyan governmeant exile. overthrown 1969.

exil princ idri libya control dissid libyan paramilitari forc origin train american intellig advisers, promis order combat col. muammar el-qaddafi, libyan leader. unit state two-year effort destabil colonel gaddafi end failur december, libyan-suppli guerrilla forc came power chad, origin 600 commando based. new chad govern ask unit state fli libyan dissid country, begin journey taken nigeria, zair final kenya. far, countri agre permanently. 400 remain commandos, disarmed, origin member libyan armi captur chad border fight 1988. volunt forc way escap p.o.w. camps. "have receiv pledg allegi leader force, princ idri step assum respons troop welfare," statement releas rome royalist libyan govern exile. overthrown 1969.

- Feed 4.4 million texts, \mathbf{w}_m , into algorithm.
- https://radimrehurek.com/gensim/models/ldamodel.html
- Approximates LDA model to back out
 - probability distribution over tokens for K topics
 - share of each topic in each text, η_m
- Use of co-occurrance to build topics.
- Reduces 0.8 million token counts to 5, 10 or 15 shares.







Aggregate Topic Shares in 4.4 Million Texts

- composition of each article m in terms of the K topics, η_m
- group of articles written in country *i* and time *t*, *M*_{*it*}
- topic shares in country i in period t is

$$\boldsymbol{\theta}_{it} = \left(\sum_{m \in M_{it}} \boldsymbol{\eta}_m \boldsymbol{N}_m + \alpha\right) / \left(\sum_{m \in M_{it}} \boldsymbol{N}_m + \boldsymbol{K}\alpha\right)$$
(1)

where $\sum_{m \in M_{it}} N_m$ is simply the total number of articles

- α enters here as the strength of the prior
- We estimate topics for each sample $T \in \{2000Q1, 2000Q2, ...2016Q4\}$ separately
- Use quarterly and yearly aggregates θ_{it} for forecasting.

"economics and trade politics" and "business and energy" in Japan



"crime and politics" and "conflict/security" in Japan



"crime and politics" and "conflict/security" in Afghanistan



both conflict topics (K = 15) in Afghanistan



• Train with data available in T

$$y_{it+1} = F(\mathbf{h}_{it}, \boldsymbol{\theta}_{it}) \tag{2}$$

where y_{it+1} is the onset of armed conflict in the quarter t+1.

- h_{it}: set of history dummies capturing post-conflict dynamics
- h_{it} also includes dummy for low-level violence and conflict in neighboring countries
- θ_{it} : share of topics
- Calculate predicted values \hat{y}_{iT+1}
- Do this for all years T = 2000Q1 2016Q4.
- To generate F(.) we feed random forests, neural network and logit regression into an ensemble

• Use the conflict history dummies as a benchmark model.

$$y_{it+1} = F(\mathbf{h}_{it}) \tag{3}$$

- How much does text add beyond this benchmark?
- Prior: text will add more in the hard problem

- We use machine learning to derive $F(\cdot)$
- This is not because it improves the fit within-sample.
- It's because we want to use the $F(\cdot)$ to produce fitted values for $\hat{y}_{iT+1}.$

- 1) we estimate a logit model
- 2) we estimate a random forest: depth 3 and 175 estimators.
- 3) we estimate a neural network: 3 layers and the number of neurons are 2, 5 and 7.
- Weight between 1), 2) and 3) is built with ensemble through soft voting.

Results: Trade-offs in Forecasting

- Need to decide on cut-off c to evaluate model:
- $\hat{y}_{iT+1} > c \rightarrow$ forecast conflict
- Two mistakes we can make:



- high cutoff c implies more false negatives
- low cutoff c implies more false positives

- ROC curves as a way to illustrate the trade-off.
- On the y-axis report the true positive rate (TPR)

$$TPR_c = rac{TP_c}{FN_c + TP_c}$$

• On the x-axis report the false positive rate (FPR)

$$FPR_c = \frac{FP_c}{FP_c + TN_c}$$

Topics help solve the hard problem



ROC curve: for any violence

Predicting onset of any violence is harder...



ROC curve: for 500+ onset

Predicting 500+ escalations is easlier.



Precision vs. True Positive Rate: 500+



- 20%: of 5 positives 1 will be a true one.
- In this literature this is really good.

Separation Plot: 500+

• Separation Plot



- machine learning like random forest is using different samples in the training
- methods are non-linear by design (example tree)

Why does it work?

- text contains some time variation
- plot predicted risk changes before onset (any violence)



Why does it work?



• plot predicted risk changes before onset (50+ conflict)

Why does it work?

• plot predicted risk changes before onset (500+ conflict)



- key is behavior of topics before onset: conflict topics go up before
- but "economics and crime", for example,



• This is even true when controlling for conflict share.

Back to ROC: 15 Topics

Onset of 1+ conflict



What Happens With Less Topics? (10)

Onset of 1+ conflict



What Happens With Less Topics (5)?

Onset of 1+ conflict



More to learn, more topics are better

Onset of 500+ conflict, 15 topics



More to learn, more topics are better

Onset of 500+ conflict, 10 topics



Onset of 500+ conflict



Compare to ICEWS

Event database of more than 6 million events



- conflict history is a great forecaster
- without conflict history forecast becomes a hard problem
- use topics to summarize massive amounts of newspaper text
- topics provide some forecast without a history
- lessons for use of unsupervised vs. supervised learning