

# Lab Notes 2

Scott Patterson

May 19, 2020

## 1 ADMIN

Hours: 9

Total: 62

Still waiting for update from purchasing department on the computer

## 2 Overview

Since the previous update, I have :

- Added UN regional groupings to the dataset. This includes dummy variable columns for AFRICA, ASIAPAC, WEOG, GRULAC, EASTEUROPE and a ‘UN\_REGION’ variable that includes each region as a separate factor.
- Generated *topical prevalence* models at k values of 18, 24, and 30 that use ‘UN\_REGION’ as a prevalence covariate.
- Used the ‘findThoughts’ function to isolate documents that best exemplify certain topics.
- Used the ‘estimateEffects’ function to investigate regional differences in the prevalence of topics relating to socioeconomic inequality. The preliminary findings suggest that Latin American states devote the most time to this topic while WEOG and Eastern European states play the topic down.
- Plotting: I had no issue generating a plot to show the general prevalence trend of inequality over time. I also plotted the difference in prevalence by region, albeit the graphic needs improvement. Visualizing regional differences over time has proven more challenging.
- Started a *topical prevalence* model using ‘year’ as an interaction term. This will allow for a deeper investigation into the relationship between time and region.

## 3 Findings

### 3.1 Covariate Addition

I had one minor issue while adding the regional covariates. After debugging several failed models, I noticed that 252 entries did not cluster into one of the defined regions. I haven’t fully investigated why, but I suspect the NA cases were from countries that broke apart or no longer exist. This may merit further investigation down the road, as we will need to deal with the addition/subtraction of new countries since 1970.

NA cases aside, the regional covariates have been added to the preprocessing stream and are now available as a covariates.

### 3.2 Prevalence Models

This section contains the prevalence model word lists for the 18 topic model. There are several topics that are potentially related to socioeconomic inequality. This analysis zooms in on Topic 5.

Based on the noted keywords, Topic 5 warrants further exploration as a potential container for tokens relevant to socioeconomic inequality. This topic also appears to capture something related to drug conflicts in Latin America. Bearing this in mind, it will be useful to see if this topic breaks down into subtopics at higher K values.

Figure 1: Topical Prevalence:  $k = 18$ ,  $\text{thresh} = 10$ ,  $\text{weighting} = \text{prob}$

(a) Topic 5 keywords: develop, trade, social, debt

Topic 1: state, nuclear, intern, peac, unit, weapon, nation, countri, peopl, world, secur, soviet, disarm, relat, arm
Topic 2: develop, nation, global, unit, countri, intern, sustain, climat, chang, secur, challeng, will, goal, commit, support
Topic 3: peopl, countri, unit, state, struggl, independ, nation, republ, govern, aggress, support, forc, world, coloni, will
Topic 4: nation, unit, intern, secur, human, organ, right, council, must, develop, will, state, member, reform, need
Topic 5: develop, countri, econom, intern, world, trade, economi, must, social, product, per, increas, debt, resourc, problem
Topic 6: countri, intern, peac, african, organ, nation, peopl, will, africa, republ, develop, communiti, state, world, govern
Topic 7: intern, peac, secur, countri, region, state, develop, peopl, effort, nation, will, unit, achiev, arab, also
Topic 8: israel, peac, state, unit, nation, arab, intern, resolut, palestinian, secur, peopl, right, aggress, territori, lebanon
Topic 9: intern, countri, nation, africa, south, develop, peopl, econom, peac, unit, continu, world, problem, namibia, state
Topic 10: nation, develop, unit, countri, will, world, intern, problem, organ, econom, assembl, session, state, general, govern
Topic 11: will, unit, must, communiti, right, nation, european, peac, problem, intern, effort, negoti, europ, year, also
Topic 12: terror, peopl, peac, will, must, nation, terrorist, secur, unit, world, human, year, govern, right, afghanistan
Topic 13: nation, develop, peac, unit, countri, world, will, intern, new, region, effort, econom, peopl, hope, cooper
Topic 14: intern, countri, peac, nation, govern, state, will, peopl, right, unit, polit, america, american, latin, respect
Topic 15: africa, peac, african, nation, govern, unit, countri, intern, communiti, develop, continu, will, conflict, secur, peopl
Topic 16: world, will, can, one, nation, peopl, must, war, time, countri, now, mani, human, year, power
Topic 17: nation, intern, unit, region, countri, state, secur, cooper, will, econom, republ, develop, organ, european, support
Topic 18: island, nation, unit, state, small, new, govern, pacif, develop, region, countri, will, peopl, continu, support

Figure 2: Topical Prevalence:  $k = 18$ , thresh = 10, weighting = frex

(a) Topic 5 keywords: debt, export, adjust, flow

Topic 1: soviet, nuclear, socialist, poland, disarm, weapon, prohibit, race, space, outer, detent, missile, armament, mongolian, treaty
Topic 2: mdgs, post-, peacebuild, bahama, ki-moon, gender, doha, saint, climat, millennium, hivaid, inclus, high-level, sustain, theme
Topic 3: imperi, imperialist, vietnames, revolutionari, nam, viet, colonialist, reactionari, victori, fascist, cuban, kampuchea, cuba, khmer, capitalist
Topic 4: peacekeep, oper, fiftieth, reform, canada, non-prolifer, prolifer, mandat, court, council, peace-keep, prevent, personnel, organization', statut
Topic 5: drug, debt, market, export, per, product, rate, cent, trade, growth, adjust, industri, traffick, trinidad, flow
Topic 6: chad, niger, mali, burundi, seneg, rwanda, zair, benin, togo, gabon, comoro, equatori, congo, burkina, cameroon
Topic 7: yemen, morocco, tunisia, kuwait, emir, iraqi, libya, arab, libyan, bahrain, kingdom, iraq, gulf, saudi, arabia
Topic 8: israel, lebanon, lebanes, zionist, isra, iranian, iran, jerusalem, jewish, palestin, arab, occup, occupi, jordan, jew
Topic 9: namibia, namibian, apartheid, pretoria, swapo, black, south, racist, front-lin, fortieth, non-align, co-oper, kampuchea, africa, indian
Topic 10: rhodesia, connexion, sea-b, waldheim, program, kurt, racial, unctad, dialog, xxv, sea, twenty-fifth, honor, favor, -develop
Topic 11: ireland, cyprus, malta, turkey, turkish, european, austria, europ, greec, mediterranean, cypriot, itali, greek, belgium, northern
Topic 12: terrorist, pakistan, terror, sri, lanka, kashmir, taliban, extremist, syria, india, timor-lest, attack, kill, muslim, afghanistan
Topic 13: japan, asean, organis, thailand, nepal, cambodia, myanmar, korea, peninsula, philippin, swaziland, indonesia, korean, people', south-east
Topic 14: panama, guatemala, bolivia, hondura, paraguay, ecuador, peru, costa, venezuela, latin, spain, rica, chile, american, argentina
Topic 15: sierra, leon, liberia, ethiopia, malawi, uganda, somalia, kenya, ghana, sudan, eritrea, somali, liberian, nigeria, gambia
Topic 16: tell, someth, hear, told, thing, truth, els, perhap, simpli, simpl, word, think, dream, man
Topic 17: azerbaijan, tajikistan, croatia, moldova, kazakhstan, belarus, georgia, herzegovina, bosnia, ukraine, turkmenistan, armenia, kosovo, latvia, macedonia
Topic 18: zealand, fiji, barbado, australia, papua, pacif, solomon, island, mauritius, iceland, fish, caledonia, maldiv, dominica, samoa

Figure 3: Topical Prevalence:  $k = 18$ , thresh = 10, weighting = lift

(a) Topic 5 keywords: hyperinfl

Topic 1: ssr, binari, ceausescu, ilyich, jaruzelski, presidium, songun, soviet-french, soviet-indian, space-bas, todor, zhivkov, byelorussian, shcherbitski, wojciech
Topic 2: oec, "effect, bangla, barbados', bioenergi, debt-gdp, fast-start, grenada', lower-middle-incom, lucia', mdg, pan-caribbean, police-contribut, unclo, erika
Topic 3: anti-imperi, battambang, cunen, glorieus, hoxha, kompong, libertaa, lon, marien, movimento, ngouabi, nol, tonkin, tse-tung, enver
Topic 4: todayâ, cost-cut, estonia', mazowiecki, secretary-generalâ, rainier, ahern, berti, burmamyanmar, â€, subsidiar, yearâ, monaco', amsterdam, marino
Topic 5: marijuana, ecuador', bradi, surinames, peru', leonel, cardoso, eclac, fujimori, quechua, garcia, narco-terror, bolivia', brazil', hyperinfl
Topic 6: niger', faso, abdoulay, ange-félix, bata, benin', biya, bongo, bozizé, cameroon', cfa, comor, compaor, compaoré, debi
Topic 7: abdrabuh, abidin, al-jab, al-saud, alija, arab-islam, bab, bahraini, emirates', emirati, hamad, mansour, moroccan-spanish, mousa, omani
Topic 8: judea, judea, samaria, zion, aqsa, dayan, irgun, kafr, khomeini, mosh, yasin, euphrat, balfour, yom, kippur
Topic 9: -payment, self-extinct, nimeiri, swapoj, traor, turnhall, xix, kountch, seyni, administrator-gener, billion—, botha, pac, pinié, concerned—
Topic 10: aaddl, addlj, non-arma, unctadj, guide-lin, addl, counter-accus, uthant, schumann, malik, s-vii, brooksrandolph, hambro, -kill, us-
Topic 11: austro-italian, famagusta, ankara, bizon, cyprus', denktash, ellemann-jensen, german-speak, rauf, turkish-cypriot, varosha, stoel, tyrol, tyrolean,
anglo-irish
Topic 12: shia, afghan-l, al-assad', aylan, benazir, drone, inter-servic, jirga, kandahar, lanka', ltte, pakistan', qaeda, loya, yazidi
Topic 13: arf, nam', nepales, uncdf, unv, yasushi, pattaya, cambodia', political-secur, marco, co-prosper, nepal', mindanao, lumbini, akashi
Topic 14: anastasio, campin, inter-ocean, itaipu, amphictyon, asuncion, banzer, belaund, belisario, betancur, caldera, debayl, humberto, interocean, montevideo
Topic 15: ajj, eritrea', ethiopia', gambia', gambian, ghana', igad, jammeh, kabbah, ketumil, leonean, mele, muluzi, museveni, nigeria'
Topic 16: néstor, sleepless, chess, smell, shakespeare, morazán, guy, lazi, millionaire, sheep, unidentifi, messiah, robot, craze, theolog
Topic 17: abkhaz, abkhazian, albania', aral, armenia-azerbaijan, armenia', ashgabat, azerbaijan', azerbaijani, baku-tbilisi-ceyhan, baku-tbilisi-kar, bulgaria', chairmanship-offic, chisinau, cica
Topic 18: drift-net, driftnet, fiji', fijian, honiara, hydrofluorocarbon, maori, marshalles, matignon, melanesian, moresbi, noumea, nouméa, palau', pohnpei

Figure 4: Topical Prevalence:  $k = 18$ , thresh = 10, weighting = score

(a) Topic 5 keywords: develop, debt

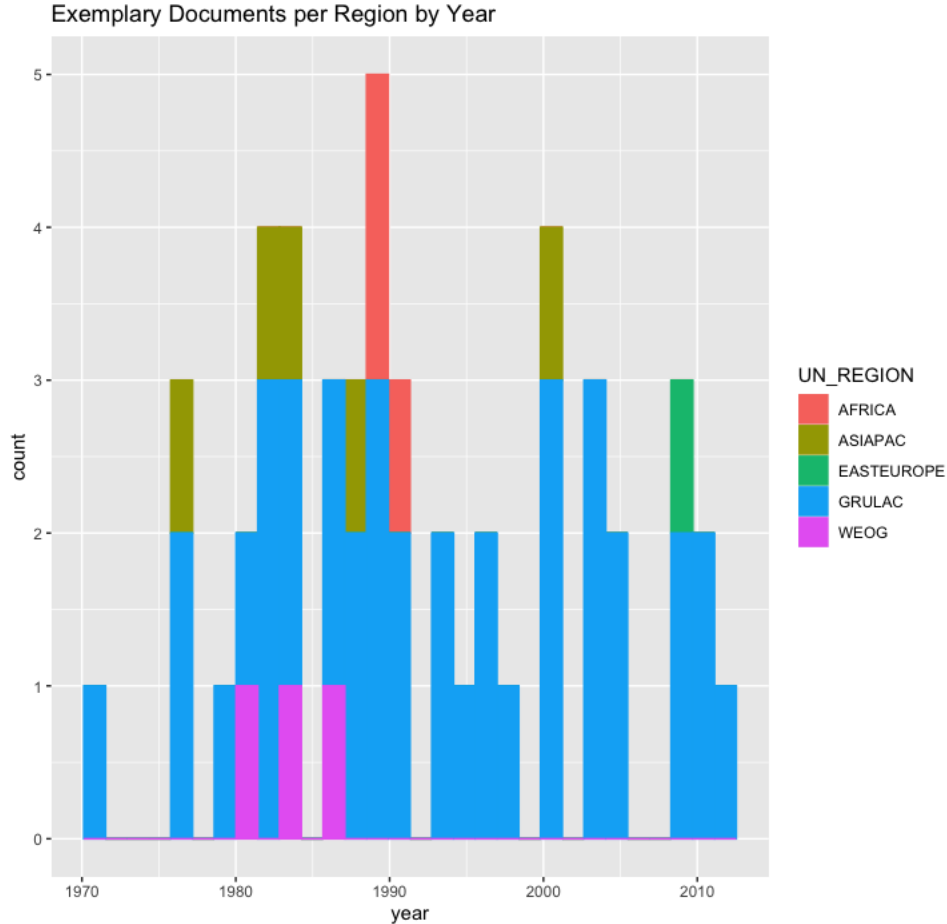
Topic 1: soviet, co-oper, mongolian, ssr, nuclear, romania, detent, poland, socialist, disarma, byelorussian, ukrainian, ussr, favor, weapon
Topic 2: mdgs, post-, millennium, hivaids, peacebuild, bahama, ki-moon, doha, global, develop, sid, gender, darfur, emiss, caricom
Topic 3: imperialist, kampuchea, imperi, viet, nam, lao, vietnames, coloni, zionist, racist, peopl, albania, viet-nam, colonialist, capitalist
Topic 4: peacekeep, millennium, unit, peace-keep, marino, reform, nation, npt, council', landmin, global, organization', annan, kofi, council
Topic 5: drug, tobago, trinidad, develop, debt, caribbean, traffick, jamaica, haiti, surinam, countri, trade, econom, market, per
Topic 6: niger, burundi, benin, chad, mali, guinea, togo, zair, equatori, rwanada, burkina, african, faso, oau, comoro
Topic 7: arab, emir, kuwait, morocco, yemen, iraqi, tunisia, islam, libyan, oman, sultan, bahrain, tunb, syrian, al-qud
Topic 8: zionist, arab, israel, islam, lebanes, lebanon, isra, palestinian, iran, iranian, moslem, palestin, jew, jewish, occup
Topic 9: co-oper, kampuchea, namibia, racist, namibian, swapo, africa, apartheid, pretoria, south, dialog, kampuchean, oau, zimbabw, contadora
Topic 10: co-oper, rhodesia, connexion, sea-b, waldheim, dialog, viet-nam, favor, kurt, detent, coloni, endeavor, honor, xxv, smith
Topic 11: co-oper, cyprus, turkish, malta, turkey, cypriot, ireland, greec, europ, european, mediterranean, csce, austria, greek, belgium
Topic 12: taliban, terrorist, terror, sri, pakistan, timor-lest, syria, lanka, afghanistan, islam, syrian, muslim, afghan, kashmir, iran'
Topic 13: nepal, thailand, asean, myanmar, lao, cambodia, swaziland, viet, japan, organism, people', nam, forty-eighth, korean, thai
Topic 14: panama, paraguay, bolivia, ecuador, peru, venezuela, hondura, costa, guatemala, rica, panamanian, american, bolivian, dominican, chile
Topic 15: malawi, sierra, liberia, leon, ecowa, liberian, somalia, african, somali, africa, lesotho, ethiopia, sudan, uganda, igad
Topic 16: world, one, can, man, war, let, know, god, even, men, want, power, say, must, ask
Topic 17: azerbaijan, croatia, tajikistan, kazakhstan, moldova, bosnia, belarus, herzegovina, turkmenistan, osc, albania, armenia, kyrgyzstan, ukraine, kosovo
Topic 18: fiji, zealand, papua, barbado, island, pacif, solomon, samoa, australia, iceland, caledonia, maldiv, guinea, mauritius, dominica

### 3.3 findThoughts

findThoughts is a function that isolates documents that contain the highest prevalence of a specified topic. The output of the function is a vector of indices that correspond to the location of the specified documents within the corpus. Substantively, this function aims to identify documents from the corpus that best exemplify the topic in question.

This figure depicts the number of exemplary documents of Topic 5 per year by region. Interestingly, 38/50 speeches with the highest proportion of Topic 5 were delivered by Latin American (GRULAC) countries. The document frequency peaked in 1990 and there were relative peaks in the early 80s and 00s. The only WEOG documents appear in the early to mid-80s.

Figure 5: Exemplary Documents per year by region



UN_REGION	Number of Exemplary Documents
AFRICA	3
ASIAPAC	5
EASTEUROPE	1
GRULAC	38
WEOG	3

The best way to act on this finding is to qualitatively analyze a group of the exemplary documents to get a sense of the themes. I can provide a list of these documents at your request.

### 3.4 estimateEffect

estimateEffect is a function used to estimate the effect of a specified covariate on the prevalence of a topic. In this case, I specified the factor variable 'UN\_REGION' plus the non-parametric 'year'. I can't yet fully explain the non-parametric "spline" that captures year in the model, but to my best understanding, this means that there is

no presumption of a linear relationship (thus non-parametric) between topic prevalence and year. Running these models made me realize that to capture the effect of year, I'll need to try it as an interaction term in a separate model.

Examining the effects table for the *magnitude* and *direction* of effects yields interesting results.

Figure 6: estimateEffect: Regional Grouping v. Topic 5 Prevalence

**Topic 5:**

**Coefficients:**

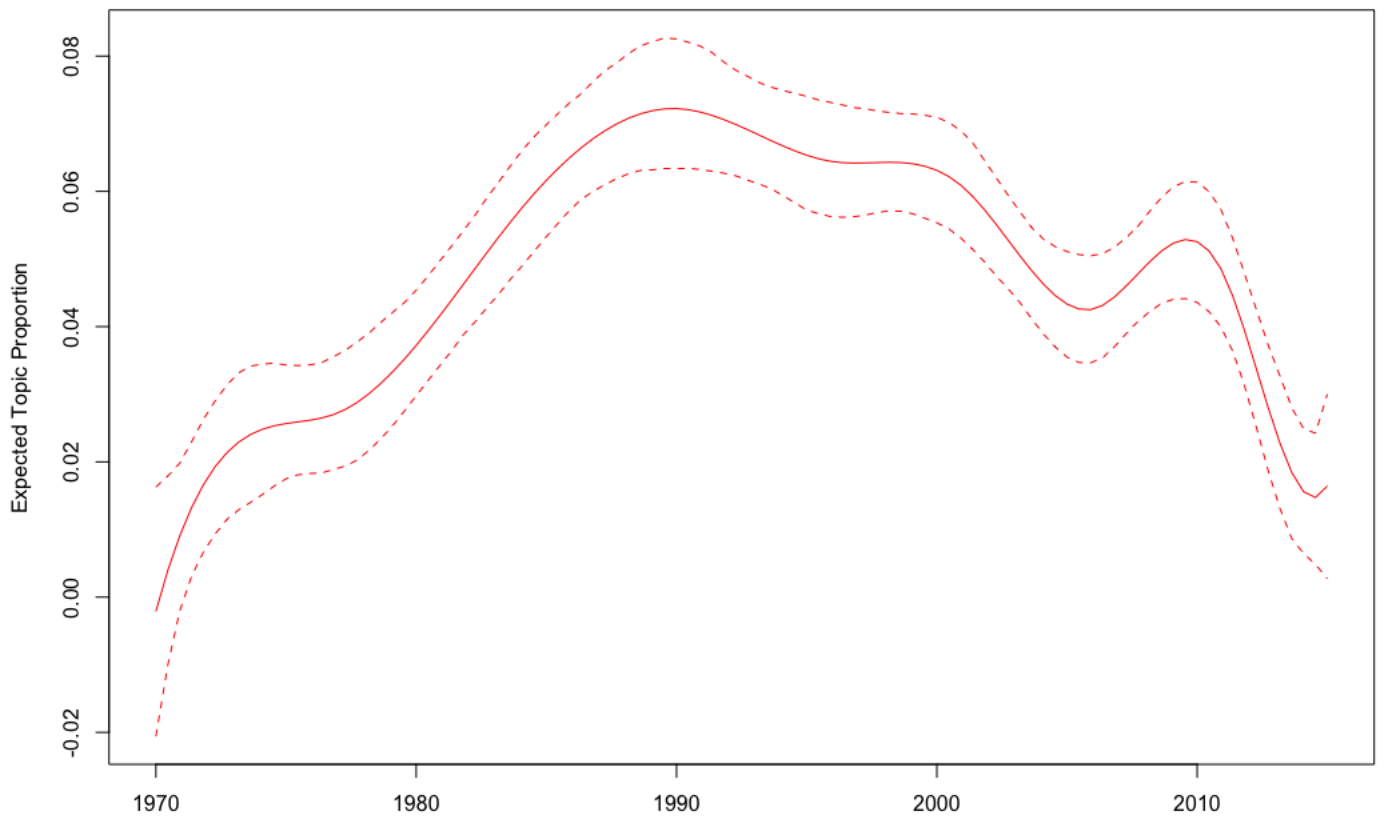
	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-0.002363	0.009419	-0.251	0.80191	
UN_REGIONASIAPAC	-0.002227	0.003090	-0.721	0.47116	
UN_REGIONEASTEUROPE	-0.022634	0.004186	-5.407	6.63e-08	***
UN_REGIONGRULAC	0.089397	0.003866	23.125	< 2e-16	***
UN_REGIONWEOG	-0.010798	0.003704	-2.916	0.00356	**
s(year)1	0.033800	0.017320	1.951	0.05104	.
s(year)2	0.018112	0.012693	1.427	0.15362	
s(year)3	0.061975	0.013369	4.636	3.62e-06	***
s(year)4	0.080586	0.012033	6.697	2.28e-11	***
s(year)5	0.062605	0.012472	5.020	5.30e-07	***
s(year)6	0.071133	0.011628	6.117	1.00e-09	***
s(year)7	0.029733	0.012557	2.368	0.01792	*
s(year)8	0.076799	0.013732	5.593	2.31e-08	***
s(year)9	0.009838	0.013112	0.750	0.45307	
s(year)10	0.018592	0.011732	1.585	0.11308	
---					
Signif. codes:	0	'***'	0.001	'**'	0.01
	'*'	0.05	'.'	0.1	' ' 1

There appears to be a negative relationship between Topic 5 prevalence and the WEOG and Eastern Europe groupings.

The relationship between GRULAC and Topic 5 prevalence appears positive and at higher magnitude than with any other grouping. Again, it appears that the Latin American states devote much more of their UNGA floor time talking about Topic 5 than any other regional grouping.

I was able to generate a few plots from the estimateEffect object. The first one depicts the general prevalence of Topic 5 over time. For this figure, the y-axis depicts the estimated proportion of the overall corpus devoted to topic 5 over time.

Figure 7: Topic 5 Prevalence over time



I also generated a plot that depicts Topic 5 prevalence by region and a plot for the prevalence of Topic 5 relative to other topics. However, I haven't yet been able to simultaneously capture regional differences and time, which leads me to think I need to come up with a new STM model that uses time as an interaction term not as a non-parametric additive.



Figure 8: Topic 5 Prevalence by Regional Grouping

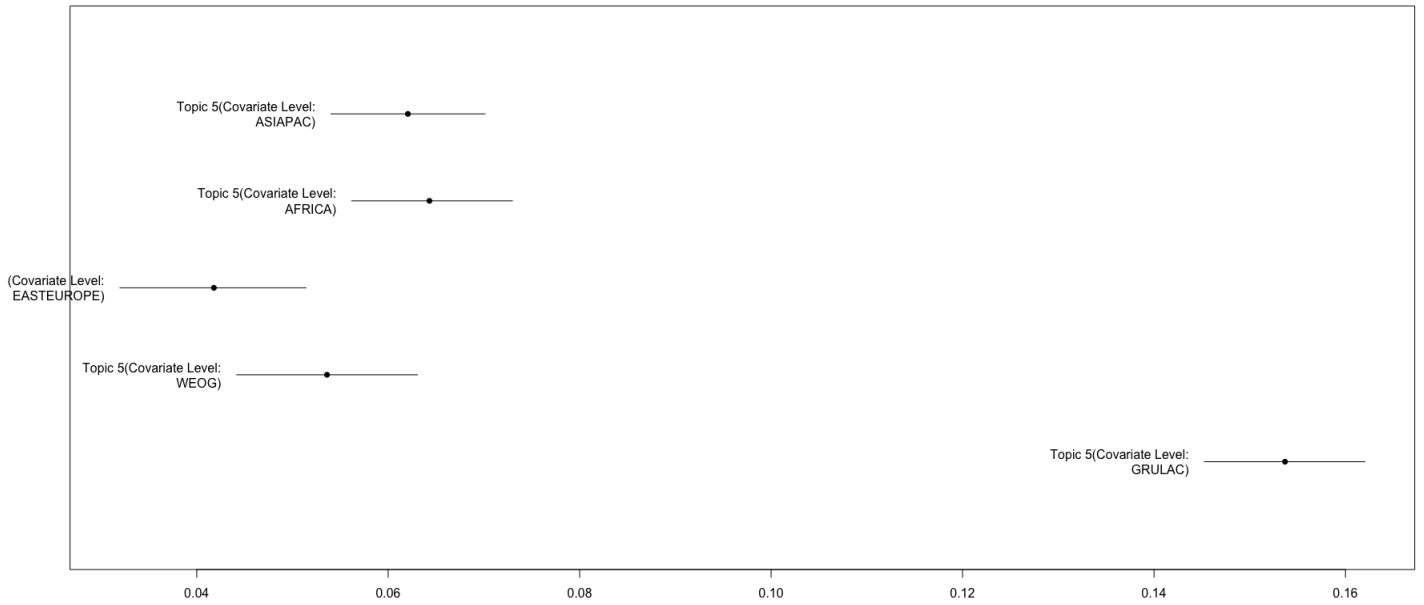
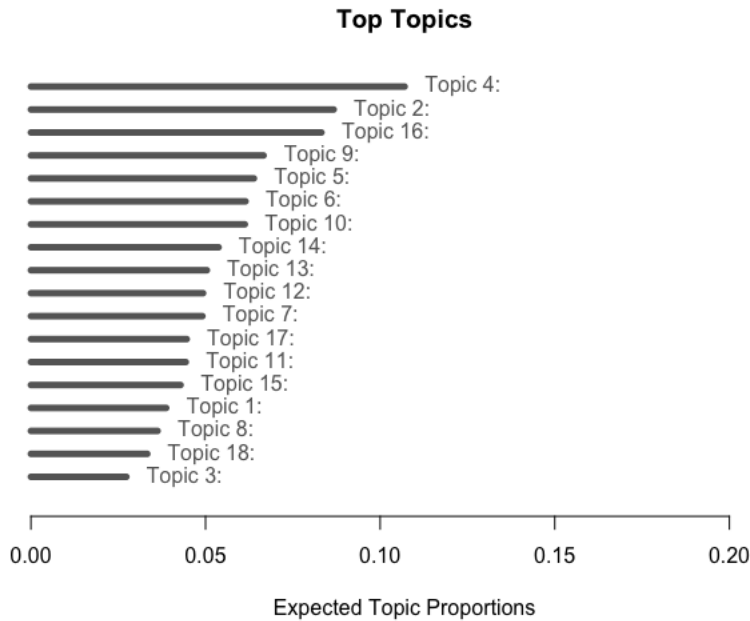


Figure 9: General Topical Prevalence

(a) Note that Topic 5 is the 5th most prevalent topic overall



## 4 Forward

My first priority for the next stretch will be to continue working on the k18 model while altering the underlying effect estimate to account for the interaction between regional groupings and time. The low hanging fruit here

will be the regional dummies, which will allow me to visualize the difference between each region and the average topical prevalence. I have already found some samples that will help me to write the syntax for this as soon as I generate the new model. The more difficult step will be to visualize the relative prevalence of each regional grouping simultaneously. This will be more difficult because the ‘estimateEffect’ object generated by STM is not easily transferable to other R packages. STM contains several basic visualization features as well as some highly advanced ones, not much in the middle. In any case, this is nothing that can’t be figured out without a bit of tinkering.

When I have finished with the step above, I’d like to replicate all of these results with the k24 and k30 models. I expect that in these models, Topic 5 will devolve into several more precise topics. This may isolate the inequality tokens from the war on drugs tokens. If not, it could also be possible that inequality and the war on drugs are rhetorically-proximate, at least in the Latin American case. We can also explore this possibility by reading through some of the documents that popped up from the ‘findThoughts’ test.

I zeroed in on Topic 5 because it seemed like a good fit at a close glance. Some of the other topics also contained tokens that may make reference to socioeconomic inequality. Now that I’ve built the pipes to run these tests, they can quickly be replicated on other topics.

The regional covariates should also be useful when returning to the *topical content* models. One issue that I’ve encountered before is that Topic 5 won’t necessarily appear in exactly the same way for a content and prevalence model. I think that I can now remedy this with an intermediate fix that I learned from Prof Erlich. When I have the new computer, I can use the extra computational power for a guaranteed remedy here, as I can generate both a content and prevalence model simultaneously.

Finally, I’m also working on a more dynamic approach to visualization that will allow you to tinker with some of the STM specifications from your web browser. I think this will be better than these snapshots that I am presenting here.

I look forward to discussing these results at your convenience.