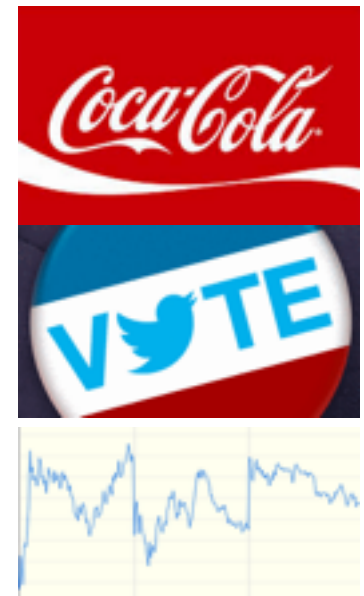


# Recall Estimation for Rare Topic Retrieval from Large Corporuses

Praveen Bommanavar (Twitter), Alek Kolcz (Twitter),  
Anand Rajaraman (Stanford)

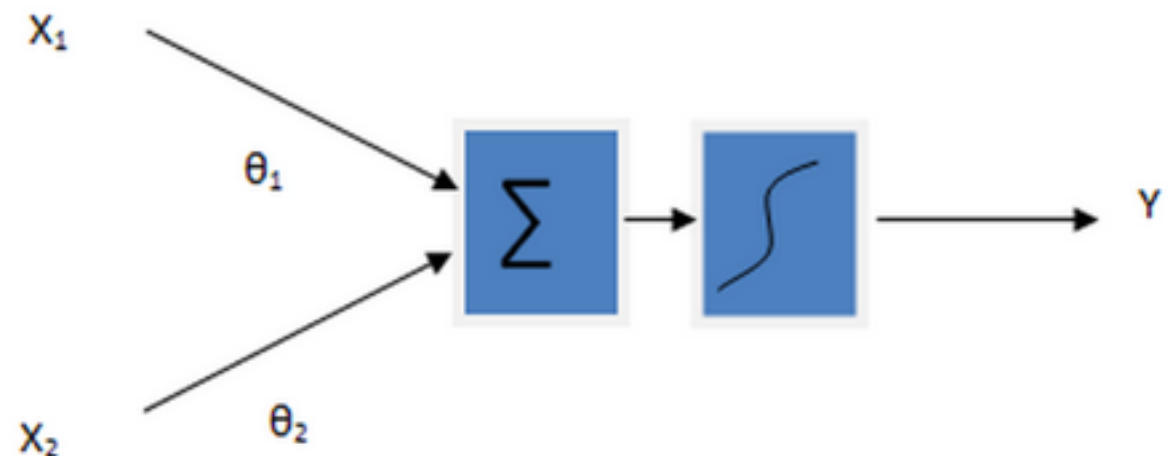
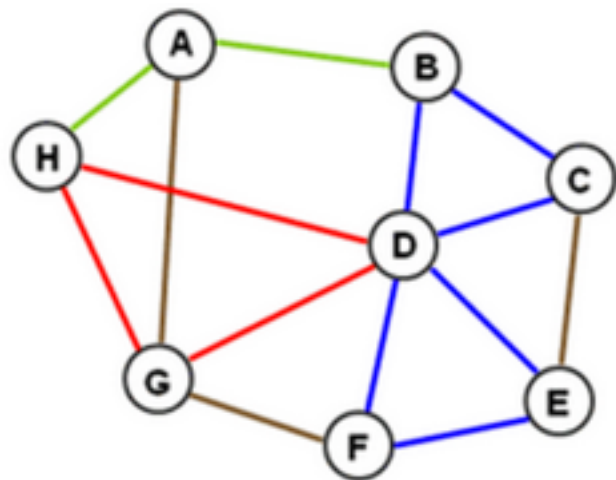
# Mining Large Corporuses

- Core offering of social media analytics companies
  - Analyze sentiment around products/brands
  - Estimate popularity of politicians
  - Uncover financial trends



# Mining Large Corpora

- Keyword filters, random walks, trained classifiers...



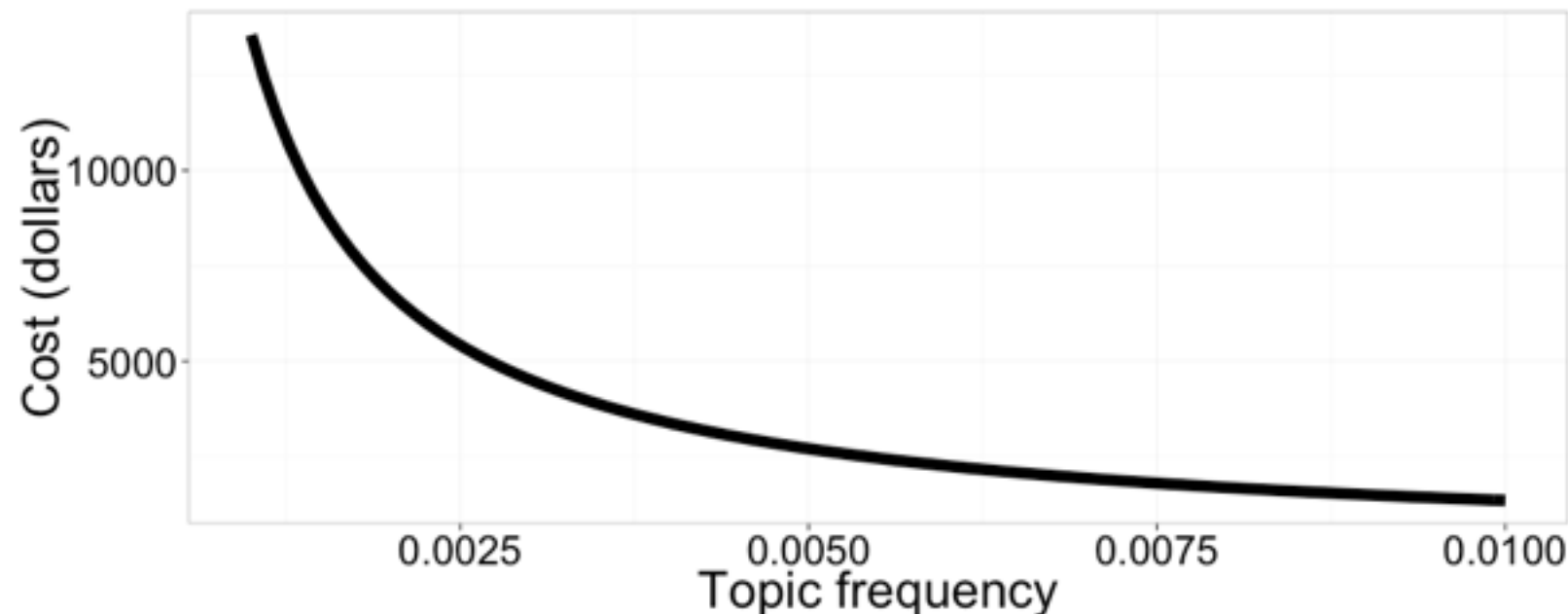
- With *any* approach: want **precision** and **recall**
  - Others: AUC, FPR, DCG, etc.

# Metrics for Rare Topics

- Precision: sample positively classified docs
  - 384 samples for 95% confidence interval of size 0.1
  - Pay approximately \$0.05 per evaluation => \$19
- Recall: sample **all** docs to find enough true positives
  - Can be very expensive if topics are rare

# Metrics for Rare Topics

- Skewed topic distribution => expensive recall est.



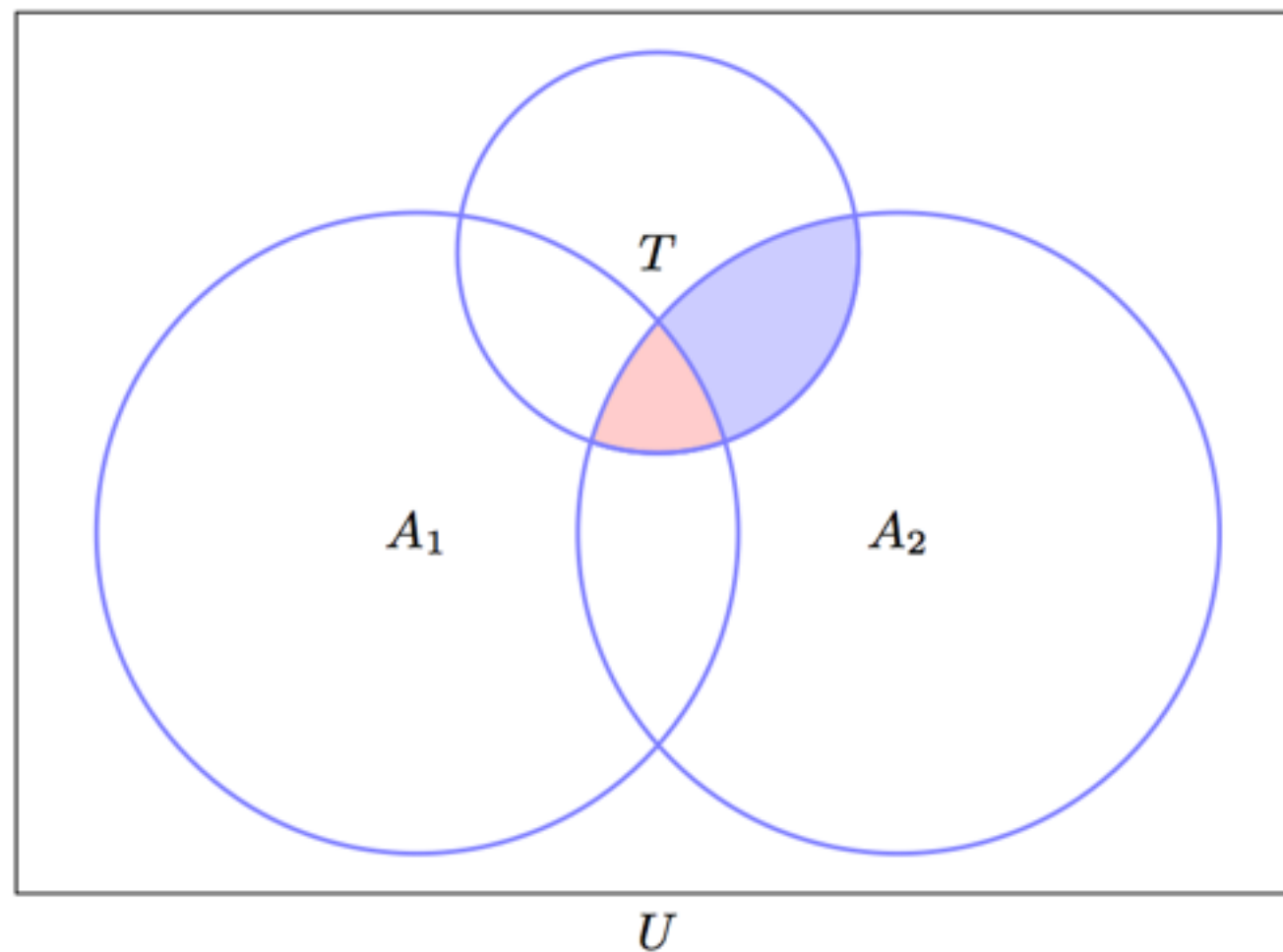
- Contribution: estimating recall on the cheap for rare topics

# Related Work

- Calibration approaches for precision [Bennett & Carvalho]
- Confidence intervals for recall (frequent classes) [Webber]
- Counting positives despite inaccurate classification (frequent classes) [Forman]
- We emphasize **cost** and **rare** classes

# Intuition

- Use pairs of *sufficiently independent* classifiers



# Conditional Independence

- $T$  = set of on topic documents
- Classifiers  $C_1, C_2$  return document sets  $A_1, A_2$

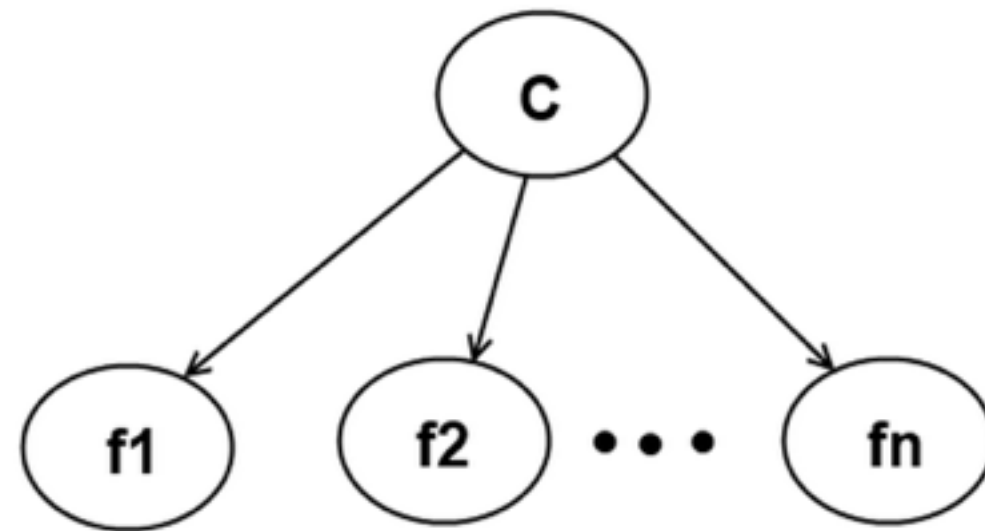
ASSUMPTION 1. (*Conditional Independence 1*) For the set of on-topic documents  $T$ ,  $C_1$  and  $C_2$  are independent classifiers. That is,

$$\delta_1 := \frac{P[A_1 \cap A_2 | T]}{P[A_1 | T]P[A_2 | T]} \approx 1.$$

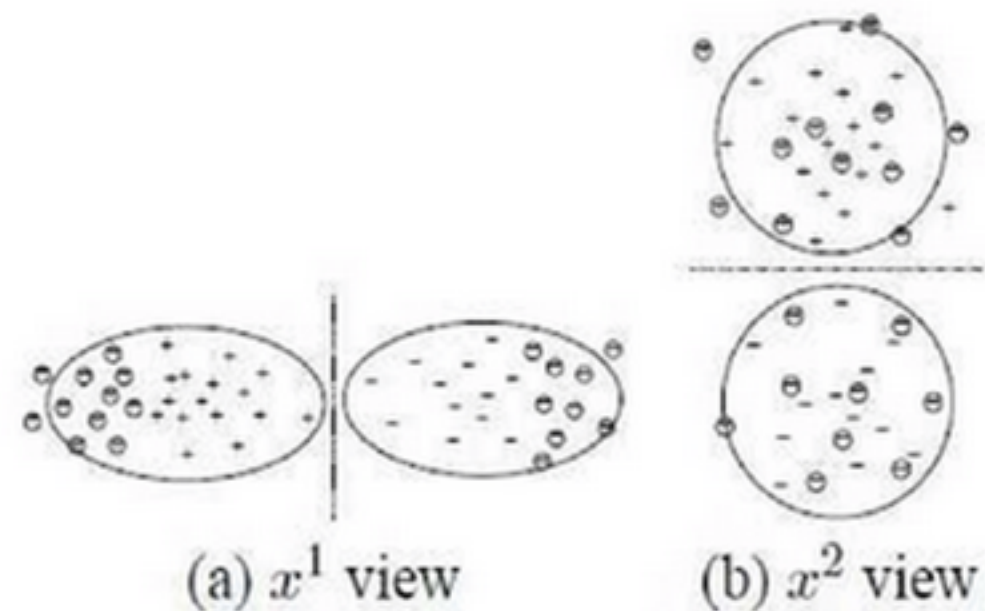


# Conditional Independence

- Naive Bayes



- Co-training



# Measuring Recall

- We can estimate recall using only precision!
- precision =  $P[T | A]$                       recall =  $P[A | T]$

$$r_1 = P[A_1|T] = \frac{P[A_1|T]P[A_2|T]}{P[A_2|T]} \approx \frac{P[A_1 \cap A_2|T]}{P[A_2|T]}$$

$$\begin{aligned} r_1 &\approx \frac{P[A_1 \cap A_2|T]P[T]}{P[A_2|T]P[T]} = \frac{P[A_1 \cap A_2 \cap T]}{P[A_2 \cap T]} \\ &= \frac{P[T|A_1 \cap A_2]P[A_1 \cap A_2]}{P[T|A_2]P[A_2]} \\ &= \frac{p_{12}|A_{12}|}{p_2|A_2|} \end{aligned}$$

# Measuring Recall cont.

- What if we don't have joint classifier precision  $p_{12}$ ?
- With a couple more assumptions, we're still in luck:

ASSUMPTION 2. (Conditional Independence 2) For the set of off-topic documents  $T^c$ ,  $C_1$  and  $C_2$  are independent classifiers. That is,

$$\delta_2 := \frac{P[A_1 \cap A_2 | T^c]}{P[A_1 | T^c]P[A_2 | T^c]} \approx 1.$$

ASSUMPTION 3. (Sparsity) The number of on-topic documents  $T$  is small, as compared with the total universe of documents  $U$ . That is,  $P[T] \ll 1$

$$r_1 \approx \frac{|A_{12}|}{p_2|A_2|} \left[ 1 - \frac{(1 - p_1)(1 - p_2)|A_1||A_2|}{|U||A_{12}|} \right]$$

# Constructing classifier pairs

- Great! Where do we get these classifier pairs from?
- Documents tend to be redundant; same info is expressed in different ways
  - Anchor text, headers, linked URLs, etc.
- Social media contains special structure

# Dataset 1: sampled Tweets

- ~1M English language Tweets from Aug 6, 2012
- topics: {apple, mars, obama, olympics, none}
- Approx \$20k budget to fully label

**Table 2: Examples of seed terms and high Jaccard similarity neighbors (Twitter statuses).**

[ <i>apple</i> , # <i>apple</i> ]	[#ios6, #ipad3, #iphone, hack, macintosh, iphones, #siri, ios, macbook, icloud, ipad, samsung, #ipodtouch, 4s, itunes, cydia, cider, #gadget, #tech, #tablet, app, connector, #mac, ...]
[ <i>mars</i> , # <i>mars</i> ]	[rover, nasa, #curiosity, #curiousity, image, mission, surface, #curiosityrover, bruno, milky, budget, @marscuriosity, gale, crater, orbiting, successfully, lands, landing, breathtaking ...]
[ <i>obama</i> , # <i>obama</i> ]	[@barackobama, barack, bush, #mitt2012, #obama2012, obamas, #dems, #gop, #military, romney, #idontsupportobama, potus, #president, administration, pres, #politics, voting ...]
[ <i>olympics</i> , # <i>olympics</i> ]	[medalist, gold, london, gb, kirani, gymnast, kate, sprinter, winning, won, #boxing, soccer, watch, watching, #usa, #teamgb, #canada, #london, javelin, nbc, match, 2012, 400m ...]

# Dataset 1 recall estimates

- All recall estimates are within 0.10 absolute error and within 15% relative error
- O(\$1000) to O(\$10)

Table 3: Experimental results: tweet keyword filters. Both recall estimation schemes are within 0.10 absolute error and 15% relative error of the true recall for all topics.

Topic	$ A $	$ A_{seed} $	$ A_{kw} $	$ A_{joint} $	$\hat{p}_{seed}$	$\hat{p}_{kw}$	$\hat{p}_{joint}$	$\hat{r}_{seed}^{(1)}$	$\hat{r}_{seed}^{(2)}$	$r_{seed}$	$\hat{r}_{kw}^{(1)}$	$\hat{r}_{kw}^{(2)}$	$r_{kw}$
Apple	3038	676	10217	420	0.655	0.247	0.774	0.129	0.166	<b>0.146</b>	0.734	0.943	<b>0.830</b>
Mars	2372	1783	7703	1433	0.904	0.264	0.938	0.661	0.704	<b>0.680</b>	0.834	0.889	<b>0.857</b>
Obama	1253	851	7400	513	0.984	0.116	0.994	0.596	0.599	<b>0.668</b>	0.609	0.613	<b>0.683</b>
Olympics	23126	4595	45705	2688	0.986	0.330	0.989	0.176	0.178	<b>0.196</b>	0.587	0.593	<b>0.653</b>



# Dataset 2: Twitter Stories

- 10.5M Discover stories from March 10, 2013:  
Tweets with hyperlinked URLs



- C1: tweet LR classifier, C2: web page LR classifier
- {ads and marketing, education, real estate and food, none}

# Dataset 2: recall estimates

- Evaluation via random sampling (prevalent enough topics)
- All recall estimates within 0.10 absolute error and most are within 15% relative error

Table 4: Experimental results: story text and webpage logistic regression classifiers. Both recall estimation schemes are within 0.10 absolute error of the true recall for all topics and most topics are within 15% relative error.

Topic	$ A_{tw} $	$ A_{web} $	$ A_{joint} $	$\hat{p}_{tw}$	$\hat{p}_{web}$	$\hat{p}_{joint}$	$\hat{r}_{tw}^{(1)}$	$\hat{r}_{tw}^{(2)}$	$r_{tw}$	$\hat{r}_{web}^{(1)}$	$\hat{r}_{web}^{(2)}$	$r_{web}$
Ads/Marketing	42073	76771	4369	0.825	0.698	0.900	0.073	0.077	<b>0.075</b>	0.113	0.120	<b>0.145</b>
Education	93292	76535	21426	0.827	0.868	0.873	0.282	0.319	<b>0.206</b>	0.242	0.275	<b>0.214</b>
Real Estate	42841	31978	12411	0.836	0.918	0.989	0.418	0.420	<b>0.413</b>	0.343	0.346	<b>0.380</b>
Food	42376	218507	20493	0.875	0.842	0.898	0.100	0.110	<b>0.122</b>	0.496	0.546	<b>0.522</b>



# Dataset 3: ODP Entries

- 110K ODP entries - similar structure to Discover

[www.criticalsoftware.com](http://www.criticalsoftware.com) - Develops and markets software products for business and mission critical information systems, and provide consulting and engineering services for enterprises.

- C1: description LR classifier, C2: web page LR classifier
- 12 topics

# Dataset 3: recall estimates

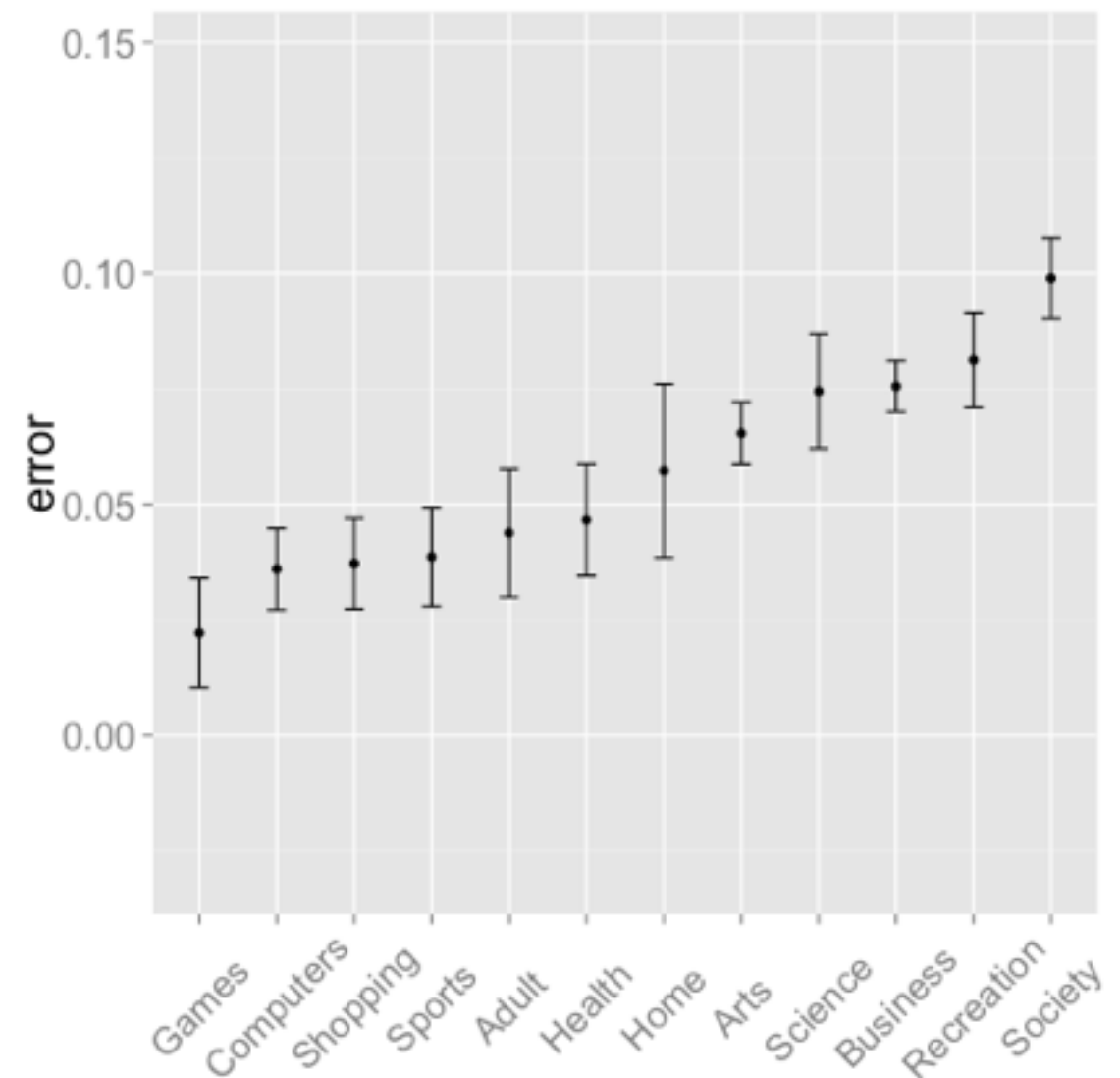
- Using joint precision directly is OK but Assumptions 2 and 3 break down

**Table 5: Experimental results: prevalence and recall estimation in ODP records. Using joint precision directly gives high fidelity recall estimates for most topics, but attempting to approximate it results in poor recall estimates.**

Topic	$ A $	$ A_{desc} $	$ A_{web} $	$ A_{joint} $	$\hat{p}_{desc}$	$\hat{p}_{web}$	$\hat{p}_{joint}$	$\hat{r}_{desc}^{(1)}$	$\hat{r}_{desc}^{(2)}$	$r_{desc}$	$\hat{r}_{web}^{(1)}$	$\hat{r}_{web}^{(2)}$	$r_{web}$
Adult	1470	10080	2815	757	0.089	0.430	0.710	0.624	0.878	<b>0.593</b>	0.839	1.180	<b>0.814</b>
Arts	13811	13719	12469	5523	0.524	0.766	0.865	0.581	0.671	<b>0.510</b>	0.771	0.891	<b>0.692</b>
Business	32304	18766	23888	10849	0.748	0.870	0.938	0.520	0.554	<b>0.443</b>	0.770	0.820	<b>0.646</b>
Computers	11235	11802	11756	4111	0.431	0.706	0.804	0.498	0.620	<b>0.453</b>	0.814	1.012	<b>0.746</b>
Games	2146	4245	3723	764	0.223	0.441	0.754	0.465	0.616	<b>0.439</b>	0.807	1.070	<b>0.766</b>
Health	5986	6180	6881	2991	0.576	0.667	0.872	0.649	0.744	<b>0.602</b>	0.837	0.960	<b>0.766</b>
Home	1546	7565	3616	643	0.118	0.299	0.574	0.594	1.035	<b>0.543</b>	0.717	1.249	<b>0.705</b>
Recreation	10846	9022	10712	4488	0.626	0.713	0.899	0.586	0.652	<b>0.505</b>	0.792	0.881	<b>0.703</b>
Science	5540	6005	7710	1706	0.380	0.462	0.658	0.474	0.720	<b>0.417</b>	0.739	1.123	<b>0.640</b>
Shopping	12386	13534	14610	3865	0.335	0.642	0.773	0.419	0.542	<b>0.375</b>	0.868	1.122	<b>0.757</b>
Society	12925	8397	11289	4071	0.627	0.720	0.922	0.501	0.543	<b>0.408</b>	0.774	0.839	<b>0.622</b>
Sports	6049	6929	6775	3197	0.550	0.708	0.910	0.664	0.729	<b>0.638</b>	0.835	0.918	<b>0.778</b>

# Dataset 3: robustness

- Estimates obtained using Assumption 1 are robust
- Random 70-30 splits



# Summary

- Have expressed recall estimates in terms of precision
- Precision is cheap to measure
- Conditionally independent classifiers can be constructed via redundancies in document structure
- **Possible future work:** Use multiple pairs of classifiers to stabilize recall estimates

# Human Evaluation

- Not exactly a turn-key system
- What could go wrong?
  - Worker impatience, fatigue & boredom, domain/lingual proficiency, laziness/scammers, definitional issues, regional differences, etc..
- What does “on-topic” even mean anyway?

# Human Evaluation

- Some remedies (not comprehensive)
  - Gold questions & agreement with other workers
  - Example answers to difficult/borderline questions (not just the easy ones)
  - Break down complex tasks into simpler ones (can't expect workers to memorize a taxonomy)
- **Communication**

# Human Evaluation

- Sometimes workers don't answer questions well, but many possible reasons. Don't simply block!
- They rate you too...

AMT Requester	Rating <a href="#">[info]</a>	Description
<b>Praveen Bommannavar</b> A1WBH67VFAHTUE HIT Group »	FAIR: 1 / 5 FAST: 1 / 5 PAY: 1 / 5 COMM: 1 / 5	Rejected my first 3 test hits within 5 minutes. He hasn't responded back yet. Sep 03 2012   <a href="#">&lt;andrewd...@h...&gt;</a>   <a href="#">flag</a>   <a href="#">comment</a>


Here's your man..an immature, unemployed university student exploiting mturk for his degree projects. If I were you, I would complain about his unethical ways to his teachers.

<https://netfiles.uiuc.edu/bommanna/www/home.htm>



# Human Evaluation

- Ran a survey about biggest pain points:
  - Communication is at the top of the list
- After some soul searching:

<b>Praveen Bommannavar</b> A1WBH67VFAHTUE HIT Group »	FAIR: 5 / 5 FAST: 5 / 5 PAY: 5 / 5 COMM: 5 / 5		Excellent requester. Wish more were like this. Have done thousands with no rejects. Clear instructions. Fast communication. Very reasonable time in taking to pay. Pay rate equals out to \$6 hr abouts if you take the time to read the tweets well and look up the occasional name to see if it is Olympic related. Probably the best requester I have worked with my 9 months or so of turking. Hope he posts regular work =] Sep 10 2012   <TurkaTRON>   <a href="#">flag</a>   <a href="#">comment</a>
<b>Praveen Bommannavar</b> A1WBH67VFAHTUE HIT Group »	FAIR: 5 / 5 FAST: 5 / 5 PAY: 5 / 5 COMM: 5 / 5		Strongly agree with all the praise here - wonderful requester. I do get the occasional rejection but since I've done about 90,000 HITs it doesn't affect me much, and I know they were genuine errors on my part. If your approval rate can take the occasional rejection, you can't do much better. To clarify, I've done about 1600 HITs and gotten 7. completely valid.



# Human Evaluation

- Email overload
  - "My dog jumped on my lap and hit my keyboard while I was working on this HIT. I'm sorry. If the answer my dog gave is wrong, I will understand the rejection. (The dog will get no treats for a week ...)"
- Other stray comments..
  - "Reading all these tweets has shattered the last little bit of hope I had for humanity. Holy hell people are stupid"