

Adaptive Language Modeling for Word Prediction

Keith Trnka

University of Delaware

advised by Kathy McCoy

Background

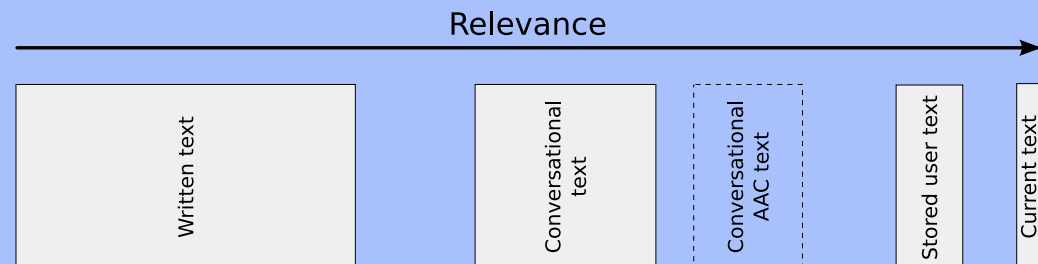
- alternative communication, slow communication rate



- word prediction speeds up communication rate
- evaluation: keystroke savings

$$KS = \frac{keys_{normal} - keys_{prediction}}{keys_{normal}} \times 100\%$$

Motivation



- ngrams sensitive to training data
- multiple uses for AAC devices
- need good (relevant) training data
- **adapt to the current text to get the most out of training data**

email excerpt

Switchboard is really low .
NNP VBZ RB JJ .

This could reflect that we chose
DT MD VB IN PRP VBD

a good corpus originally , maybe that
DT JJ NN RB , RB IN

the cleanup was more consistent
DT NN VBD RBR JJ

(I do n't think it 's any
-LRB- PRP VBP RB VB PRP VBZ DT

more advanced than the others ,
RBR JJ IN DT NNS ,

but I think I spent far more time on it
CC PRP VBP PRP VBD RB JJR NN IN PRP

paper excerpt

The self-test analysis is affected
DT JJS NN VBZ VBN

by both the size of the corpus
IN DT DT NN IN DT NN

as well as the diversity of the corpus
IN RB IN DT NN IN DT NN

, which explains the trend with Switchboard
, WDT VBZ DT NN IN NNP

: participants in the corpus collection
: NNS IN DT NN NN

were restricted to one of roughly 70 topics
VBD VBN TO CD IN RB CD NNS

, most of which are represented
, JJS IN WDT VBP VBN

in every set of Switchboard .
IN DT NN IN NNP

Adapting to Match the Topic

$$P(w | h) = \sum_{t \in \text{topics}} P(t | h) * P(w | h, t)$$

Topic Granularity

- granularity of topic labels: the size of topics; specific or general topics
- medium-grained: human-annotated, typical clusters (e.g., clothing, weather, jobs)
- fine-grained: document as topic, IR-like modeling (e.g., seasonal clothing at work)
- coarse-grained: corpus as topic, very high-level (e.g., news, chit-chat)
- evaluation (with domain variations)

	In-domain	Out-of-domain	Mixed-domain
Trigram baseline	60.35%	53.88%	59.80%
SVB (medium)	61.48% (+1.12%)	-	-
Documents (fine)	61.42% (+1.07%)	54.90% (+1.02%)	61.17% (+1.37%)
Corpora (coarse)	-	52.63% (-1.25%)	60.62% (+0.82%)

Topic Identification

- current document representation: frequency, recency, inverse topic freq.
- similarity scores: cosine (best), Jacquard, Naïve Bayes (worst)
- polarizing the scores for more discrimination
$$sim'(t, h) = \frac{sim(t, h) - \min_{t'}(sim(t', h))}{\max_{t'}(sim(t', h)) - \min_{t'}(sim(t', h))}$$
- smoothing to prevent non-zero scores for sparse topics
$$sim'(t, h) = \frac{sim(t, h) + \gamma * \min_{t'}(sim(t', h))}{\max_{t'}(sim(t', h)) + \gamma * \min_{t'}(sim(t', h))}$$
- stemming helps with sparse topics (+0.2%) but hurts for normal topics (-0.1-0.2%)

Topic Application

- using trigrams
- smooth/backoff after interpolation - interpolating frequencies
- rescaling the frequency distribution for smoothing (+0.2-0.4%)
$$\sum_w f'_{topic}(w | h) = \alpha * \sum_w f_{topic}(w | h) = \sum_{t \in \text{topics}} \sum_w f(w | h, t)$$
- binning frequencies for smoothing
- smoothing extremely sparse conditional distributions on-demand
$$\frac{f(w | h)}{f(w | h) + \lambda} \times \frac{f(w | h)}{f(h)}$$
- modeling h and t independently
$$P_{hybrid}(w | h) = P(w | w_{-2}, w_{-1}) * \left(\sum_{t \in \text{topics}} P(t | h) * P(w | t) \right)^\alpha$$

Future: style adaptation

- POS tags and pairs across styles

POS unigrams										POS bigrams									
Email					Papers					Email					Papers				
POS	f	p	d	p1/p2	POS	f	p	d	p2/p1	tags	f	p	d	p1/p2	tags	f	p	d	p2/p1
Frequent and different POS tags										Frequent and different POS bigrams									
PRP	699	0.0797399041752224	1.28	4.5649	RB	598	0.0385458295732886	0.52	0.5846	NN IN	2426	0.0766666103582	0.50	0.5788	NN IN	7406	0.0766666103582	0.50	0.5788
RB	578	0.0659365731234314	0.52	1.7106	VB	480	0.0309397963130076	0.56	0.5627	TO VB	227	0.0258955053616245	0.38	1.3870	NN NN	6190	0.03898445619827	0.61	2.0101
VB	482	0.054985169974903	0.56	1.7772	VBN	478	0.0308108804950367	0.61	1.8756	NN NN	1740	0.0198494182067077	0.40	0.4975	IN NN	4090	0.0263632847750419	0.51	1.8055
VBP	301	0.0343372119552818	0.50	1.6699	VBD	196	0.020562072966353	0.50	0.5988	JJ NN	1730	0.0197353410905772	0.40	0.5316	NN VBZ	2700	0.0174036354260668	0.49	2.0616
VBD	202	0.0230435774583619	0.58	1.8240	PRP	271	0.0174680933350522	0.28	0.2191	IN PRP	1490	0.016973003034451	0.38	1.7945	NN IN	2360	0.0152120665205621	0.48	1.8521
VBN	144	0.016427104727926	0.61	0.5332	VBD	196	0.0216337501611448	0.58	0.5483	IN NN	1280	0.0146018708647045	0.30	0.5539	NN NNS	2300	0.0148253190666495	0.38	2.9536
-LRB-	98	0.0111795573807894	0.74	2.1680	:	90	0.00580121180868893	0.61	0.5353	PRP MD	1210	0.013803331051791	0.38	1.7588	TO VB	1950	0.01269292521593	0.38	0.5509
:	95	0.0108373260323979	0.61	1.8681	-RRB-	80	0.0051566327188346	0.74	0.4613	PRP VBD	1110	0.012662598804046	0.38	1.4926	VBN IN	1860	0.0119891710712004	0.36	2.2361
-LRB-	88	0.0100387862194844	0.67	1.9967	-LRB-	78	0.00502771690086374	0.67	0.5008	VB IN	99	0.011296344969199	0.38	1.9620	IN JJ	1750	0.0112801340724807	0.62	2.2473
RP	60	0.00684462696783025	0.68	1.1683	RBR	60	0.00386747453912595	0.66	1.9943	RB VB	99	0.0050311995192661	0.38	1.9260	IN	1500	0.00966868634781488	0.38	2.5684
PRPS	56	0.00638831850330824	0.57	1.8020	PRPS	55	0.00354518499419879	0.57	0.5549	PRP	99	0.0050311995192661	0.38	1.9260	IN	1500	0.00966868634781488	0.38	2.5684
SYM	49	0.00558977869039471	0.28	2.9066	FW	31	0.00199819517854841	0.55	1.7516	PRP VBZ	99	0.0050311995192661	0.38	1.9260	IN	1500	0.00966868634781488	0.38	2.5684
WP	30	0.00342231348391513	0.12	3.5396	EX	16	0.00103132654376692	0.08	0.3117	NN VBZ	99	0.0050311995192661	0.38	1.9260	IN	1500	0.00966868634781488	0.38	2.5684
EX	29	0.00330823636778462	0.08	3.2077	WP	15	0.000966868634781488	0.12	0.2825	RB JJ	99	0.0050311995192661	0.38	1.9260	IN	1500	0.00966868634781488	0.38	2.5684
RBR	17	0.00193931097421857	0.66	0.5014	RP	13	0.000837952816810623	0.26	0.1224	PRP	99	0.0050311995192661	0.38	1.9260	IN	1500	0.00966868634781488	0.38	2.5684
FW	10	0.00114077116130504	0.55	0.5709	SYM	3	0.000193373726956298	0.08	0.0346	NN IN	99	0.0050311995192661	0.38	1.9260	IN	1500	0.00966868634781488	0.38	2.5684

- style granularity: treat each corpus as a style (e.g., Switchboard, Micase)
- style identifications: cosine similarity of POS tags and pairs
- style application: condition transition probabilities of POS ngram model

$$P_{style}(w | h) = \sum_{s \in \text{styles}} P(s | h) * \sum_{tag \in POS(w)} P(tag | tag_{-1}, tag_{-2}, s) * P(w | tag)$$