## Information Ordering

LING 573 — NLP Systems & Applications 4/30/2018

WASHINGTON





# Using External Python Packages

- If you want to use additional python packages on patas:
  - python3 -m venv \$SOME\_PATH
  - create a requirements.txt file for the packages you use in your root dir
  - use **\$SOME\_PATH/bin/python3** to run code
    - ...you can also add \$SOME\_PATH/bin to your \$PATH variable





MY PYTHON ENVIRONMENT HAS BECOME SO DEGRADED THAT MY LAPTOP HAS BEEN DECLARED A SUPERFUND SITE.



# Using External Python Packages

- This will allow:
  - local installation of packages from PyPy
  - Other user to quickly install dependencies:
    - pip install -r requirements.txt
- Standard way of supporting dependencies in python packaging







### Semi-Supervised Document Clustering

0

#### KMeans Clustering on TAC Categories 5 Clusters, 20 Runs



VASHINGTON

= tf•idf via sklearn



0

# Supervised Samples Seeding Cluster

12

= GloVe via spacy



•

20

18

16

# Semi-Supervised Document Clustering

- Code on github:
  - <u>github.com/rgeorgi/tac-clusters</u>





## Information Ordering: Combining Experts





### Integrating Ordering Preferences Bollegala et al, 2012

- Key idea:
  - Information ordering is a combination of factors.
  - Consider multiple "experts" that model different factors
  - Combine in a linear combination to determine ordering







- I. Chronological Expert
- 2. Probabilistic Expert
- 3. Topical-closeness Expert
- 4. Precedence Expert
- 5. Succession Expert



Training "Experts"



### **Basic Framework**

- Build one expert for each preference
  - Iterate through pairs of sentences (u, v) and partial summary Q
    - prefer u before v if score > 0.5
    - prefer v before u if score < 0.5
- Learn weights for linear combination
- Use greedy algorithm to produce final order



 $PREF_{total}(u,v,Q) = \sum w_e PREF_e(u,v,Q)$  $e \in E$ 





# Chronological Expert

- If sentences from two different documents with different times
  - Order by document timestamp
- If sentences from same document
  - Order by order within document
- Otherwise, no preference

WASHINGTON

 $PREF_{chro}(u,v,Q) =$ 0.5 ()

$$T(u) < T(v)$$
$$\left[D(u) = D(v)\right] \land \left[N(u) < N(v)\right]$$
$$\left[T(u) = T(v)\right] \land \left[D(u) \neq D(v)\right]$$

otherwise



## Probabilistic Expert

- Based on Lapata (2003)
- Model the probability of u preceding v in summary features:
  - POS tags
  - Dependency Structures
  - Lemmas (Smoothed)
- 0.5 is returned for equally likely outcomes, u preferred if > 0.5

$$PREF_{prob}(u,v) =$$



 $1+P(v \mid u)-P(u \mid v)$ 



# Topical-closeness Expert

- Same motivation as <u>Barzilay, 2002</u> (Clustering sentences into "themes")
   I. The earthquake crushed cars, damaged hundreds of houses, and terrified people for
  - I. The earthquake crushed cars, damaged hundreds of kilometers around.
  - 2. A major earthquake measuring 7.7 on the Richter scale rocked north Chile Wednesday3. Authorities said two women, one aged 88 and the other 54, died when the crushed
  - 3. Authorities said two women, one aged a under the collapsing walls
- | and 3 discuss theme of impact
- 2 describes magnitude and location
- Better order (2), (1, 3)





COMPUTATIONAL LINGUISTICS

# **Topical-closeness Expert**

- Q = sentences ordered thus far,  $q \in Q$
- Look at candidate sentences u, v
  - Pick one with closest similarity to already ordered sentence q
  - 0.5 if similarity is identical

WASHINGTON

 $PREF_{topic}(u,v,Q) = -$ 

 $topic(l) = \max_{q \in Q} sim(l,q)$ 

$$\begin{bmatrix} Q = \varnothing \end{bmatrix} \lor \begin{bmatrix} topic(u) = topic(v) \end{bmatrix}$$
$$\begin{bmatrix} Q \neq \varnothing \end{bmatrix} \lor \begin{bmatrix} topic(u) > topic(v) \end{bmatrix}$$

otherwise



# **Topical-closeness Expert**



- Bollegala et al use cosine similarity
  - Could use any similarity measure, suggest WordNet as an alternative



 $topic(l) = \max_{q \in Q} sim(l,q)$ 





# Precedence Expert

- With following example:
- (a) Honduran death estimates grew from 32 to 231 in the first two days, to 6,076 with 4,621 missing. (b) Honduras braced as category 5 Hurricane Mitch approached.
- (c) The EU approved 6.4 million in aid to Mitch's victims.
- (b) introduces event that is needed to understand (a), (c)
  - (a) and (c) contain information preceded by (b)











# Preceden $pre(l) = \frac{1}{|Q|} \sum_{q \in Q} \max_{p \in P_q} sim(p, l)$

• For candidate sentence *l*:

- For every already ordered sentence q:
  - Find max similarity of any sentence *p* preceding *q* in *q*'s original document to *l*
- Average this for all q.

#### • Idea:

WASHINGTON

 Sentences with maximum similarity to sentences preceding those in Q should (similarly) come first.







- No preference if there are no sentences already in Q
- If precedence of u more than v, prefer u
- Otherwise prefer *v*



# Precedence Expert





- Inverse of precedence
  - Calculate similarity of candidate with information that succeeds Q in original docs

 $PREF_{succ}(u,v,Q) = \begin{cases} 0. \\ 1 \end{cases}$ 





0.5 
$$[Q = \varnothing] \lor [succ(u) = succ(v)]$$
  
1  $[Q \neq \varnothing] \land [succ(u) > succ(v)]$   
0 otherwise

otherwise





#### **Precedence Expert**



#### **Succession Expert**





# Learning Algorithm

- Use the same algorithm to find optimal weights as **Barzilay (2002)** 
  - Namely, <u>Cohen et. al (1999)</u>
- Use model summaries to train
  - Learn optimal weights for experts given training data









# Learned Weights

#### • Optimal learned weights:



Expert	Weight
Chronological	0.327947
Probabilistic	0.000039
Topical-closeness	0.016287
Precedent	0.196562
Succession	0.444102





Correlated scores of various approaches against human judgments

#### Method

Random Ordering

**Probabilistic Ordering** 

Chronological Orderin

Proposed Method

- Probabilistic ordering is rubbish
- Chronological actually does pretty well
- Combined model with learned weights better than Chronological alone

WASHINGTON

### Results

		Spearman
	(RO)	-0.267
	(PO)	0.062
ng	(CO)	0.774
	(LO)	0.783





### Observations

- Nice ideas:
  - Combines multiple sources of ordering preferences
  - Weight-based integration
- Issues:
  - Sparseness everywhere
    - Ubiquitous word-level cosine similarity
    - Probabilistic models
  - Score handling











Entity-Based Ordering





### Entity-Based Ordering Barzilay & Lapata (2005, 2008)

- Continuing to talk about same thing(s) lends cohesion to discourse
- Incorporated variously in discourse models
  - Lexical chains: Link mentions across sentences
    - Fewer lexical chain crossings → fewer shifts in topics
  - Salience hierarchies, information structure
    - Subject > Object > Indirect > Oblique ...
  - Centering model of Coreference
    - Combine grammatical role preference with
    - Preference for types of references/focus transitions







### Entity-Based Ordering Barzilay & Lapata (2005, 2008)

### • Idea:

- Leverage patterns of entity (re)mentions
- Intuition:
  - Capture local relations between sentences, entities
  - Model cohesion of evolving story







### Entity-Based Ordering Barzilay & Lapata (2005, 2008)

#### • Pros:

- Largely delexicalized
  - Less sensitive to domain/topic than other models
- Can exploit state-of-the-art syntax, coreference tools







### Entity Grid oss sentences of:

- Compact representation across sentences of:
  - Mentions
  - Grammatical Roles
  - Transitions







#### The Entity Grid Jovernment **Jepartmen** Microsoft Products Brands Case Netscape Software Tactics Competit Markets Evidence [rial 2 - - s o - - - s o o - - - 33

6

- Rows = sentences
- Columns = discourse entities
- Values = grammatical role of mention in sentence
  - (S)ubject
  - (O)bject
  - X (other)

WASHINGTON

- (no mention)
- Multiple mentions Take highest grammatical ranking (S > O > X)

PROFESSIONAL MASTER'S IN COMPUTATIONAL LINGUISTICS



Suit



	Department	Trial	Microsoft	Evidence	Competitors	Markets	Products	Brands	Case	Netscape	Software	Tactics	Government	Suit	Earnings	
1	S	0	S	X	0	_	_		_	_	_	_	_		_	1
2	—	_	0	_	—	X	S	0	_	—	_	_	_	_	—	2
3	_	_	S	0	_		_	_	S	0	0	_	_	_	_	3
4	—	_	S	—	—		_	_	_	_	_	S	_	_	_	4
5	_	_	_	_	_		_	_	_	_	_	_	S	0	_	5
6	_	X	S	_	_	_	_	_	_	_	_	_	_	_	0	6

- 2 competitive enough to unseat [established brands]o.
- 3 merging [browser software]o.
- 4 [Microsoft]s claims [its tactics]s are commonplace and good economically
- 5 [collusion] x is [a violation of the Sherman Act] o.
- [Microsoft]s continues to show [increased earnings]o despite [the trial]x. 6

WASHINGTON

[The Justice Department]s is conducting an [anti-trust trial]o against [Microsoft Corp.]x with [evidence]x that [the company]s is increasingly attempting to crush [competitors]o

[Microsoft]o is accused of trying to forcefully buy into [markets]x where [its own products]s are not

[The case]s revolves around [evidence]o of [Microsoft]s aggressively pressuring [Netscape]o into

[The government]s may file [a civil suit]o ruling that [conspiracy]s to curb [competition]o through



#### Intuitions

- Some columns dense: focus of text (e
  - Likely to take certain roles, e.g. S, O
- Others **sparse**: likely other roles(X)
- Local transitions reflect structure, topic
- local entity transitions:  $\{S,O,X,-\}^n$ 
  - Continuous column subsequences ("role n-grams"?)
  - Compute probability of sequence over grid:
  - # of occurrences of that type/# of occurrences of that length





e.g. Microsoft)		Department	Trial	Microsoft	Evidence	Competitors	Markets	Products	Brands	Case	Netscape	Software	Tactics	Government	Suit
	1	S	0	S	X	0	_	-	_			_	_	_	_
	2	_	-	0	-	-	X	S	0	_	_	_	-	_	_
	3	_	-	S	0	-	_	-	_	S	0	0	-	_	_
c shifts	4		-	S	-	-	_	-	_		_	_	S	_	_
	5	_	-		-	-	_	-	_		_	_	_	S	0
	6	_	X	S	-	-	-	-	_	_	_	—	_	_	_





#### Vector Representation 0 0 0 .03 0 0 0 $d_1$ 0 0 0 .02 0 .07 0 $d_2$ .02 0 0 .03 0 0 0 $d_{3}$

• Document vector:

- Length = # of transition types
- Values = Probabilities of each transition type
- Can vary by transition types
  - e.g. most frequent; all transition of some length, etc.



	$\mathbf{v}$	0	$\mathbf{X}$		$\mathbf{V}$	$\bigcirc$		I	
0	X	X	X	X	I	I		I	
.02	.07	0	0	.12	.02	.02	.05	.25	
.02	0	0	.06	.04	0	0	0	.36	
.06	0	0	0	.05	.03	.07	.07	.29	



