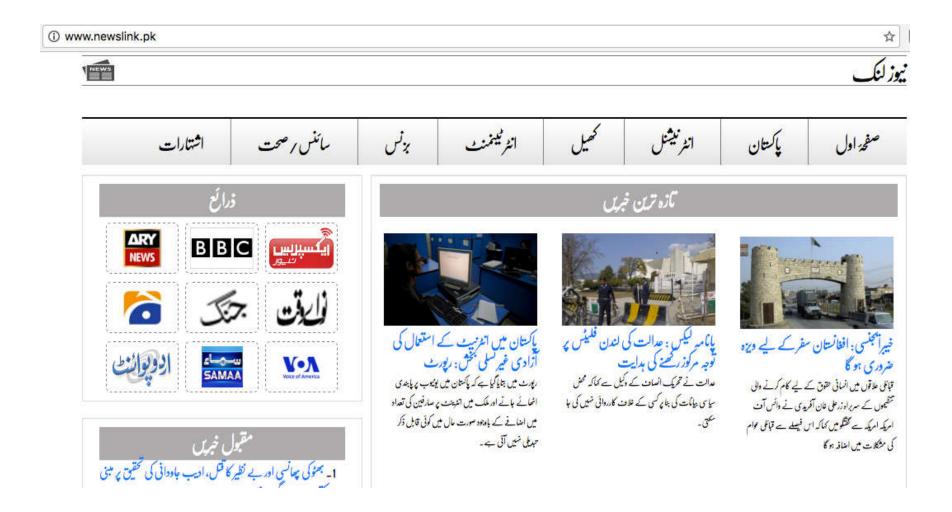
Clustering Urdu News Using Headlines

Samia Khaliq, **Waheed Iqbal**, Faisal Bukhari, Kamran Malik PUCIT, University of the Punjab, Lahore, Pakistan

Motivation



Motivation (Cont.)



Introduction

- We are actively scraping Urdu news from news agencies including Jang, Geo, Express, Nawa-e-Waqat etc.
- Currently we have around 0.5 million news scrapped from different news websites.
- Our aim is to develop a portal that can provide Urdu news analytics.
- In this work, we presented our algorithm to automatically cluster similar Urdu news.

News Scrapping

- We use Selenium to scrap the news websites.
- We configured a cron job to run the scrapper after every 4 hours.
- We scrap around 700 to 1200 news every day.

Pre-Processing

News headlines are cleaned in this step by removing stop words and identifying tokens.

News	Stop Words	Tokens
حکومت جوڈیش کمیین		· · · · ·
کے قیام پر رصنا مند ہے، عارف علوی	4.2.4	حکومت، جوڈیشل، کمیثن، قیام، رحنامند، عارف، علوی
وزیر اعلی پنجاب سے وزیر داخلہ چوہدری نثار علی خان		
کی ملاقات نان کی ملاقات	ہے،کی	وزیر، اعلیٰ، پنجاب، وزیر، داخلہ، چوہدری، نثار، علی، خان، ملاقات
سعید اجل کا ورلڈ کپ میں حصہ یہ لینے کا فیصلہ	کا، میں، یذ، لینے	سعيد، اجل، ورلدُ کپ، حصه، فيصله
مہمندا بجنسی، بینظیر بھٹو شہید کی ساتوں بر سی عقیدت واحترام سے منائی گئی	کی، و، ہے، گھتی	مهمند، ایجنسی،، بینظیر، بھٹو، شہید، ساتویں، بر سی، عقیدت، اخترام، منائی
ایم کیوایم نے نیشنل کاونٹر ٹیریزم پالیسی پر سفارشات کو حتمی شکل دیدی	نے، پر، کو، دیدی	ایم، کیو ایم، نے نیشنل، کاونٹر، نیریزم، پالیسی، سفارشات، حتمی، شکل
ملیشیاکی نجی ہوائی کمپنی کا طیارہ لاپتہ ہو گیا	کی، کا، ہوگیا	مليشي <mark>ا،</mark> نجی، ہوائی، کمپنی، طیارہ، لاپن <i>ة</i>
میلبرن ٹیسٹ؛ کوبلی اور رہانے ڈٹ گئے	اور، گخ	میلبرن، نییٹ:، کوہلی، رہانے، ڈٹ

New Similarity

Given two news *i* and *j* we need to identify similarity between them:

$$st_{avg} = \frac{st_i + st_j}{2}$$

$$S_{i,j} = \frac{m_{i,j}}{st_{avg}}$$

 st_i - size of tokens list of news i

 st_j - size of tokens list of news j

 st_{avg} - average size of token lists of news i and j $m_{i,j}$ - count of similar tokens of news i and news j

	gorithm 1: getRelatedNewsList
I	nput: n_t , news headline which is a source headline and we need to identify list of similar news
	to this news from our dataset
I	Result: Algorithm return a list of news similar to given news n_i
1 1	egin
2	$relatedNewsList \leftarrow null$
3	$t \leftarrow 0.5$
4	$datasetList \leftarrow getDatasetList()$
5	for $j = 0$ to datasetList.size do
6	$n_1 \leftarrow datasetList[j]$
7	$S_{i,j} \leftarrow getSimilarityScore(n_i, n_j)$
8	if $S_{t,t} \ge t$ then
9	$relatedNewsList.add(n_j)$
10	end
11	end
12	return relatedNewsList
13 e	
	$etSimilarityScore(n_i, n_j)$
15 b	egin
16	$tl_i = _getTokensList(n_i)$
17	$tl_j = _getTokensList(n_j)$
18	$st_i = tl_i.size$
19	$st_j = tl_j.size$
20	$S_{i,j} \leftarrow 0$
21	$m_{i,j} \leftarrow 0$
22	for $x = 0$ to st_i do
23	for $y = 0$ to st_j do
24	if $tl_t[x]$ matches $tl_j[y]$ then
25	$m_{i,j} \leftarrow m_{i,j} + 1$
26	end
27	end
28	end
29	if $m_{i,j} > 0$ then
30	$avg = (st_i + st_j)/2$
31	$S_{i,j} = m_{i,j}/avg$
32	end
33	return $S_{i,j}$
34	end

Ground Truth

- We randomly selected 100 news from each of the following categories:
 - \circ International
 - \circ National
 - \circ Health
 - o Business
 - o Entertainment
- Each news category is given to three different persons who marked related news manually. Then we identified the common clusters as ground truth in each of the category.

Evaluation Criteria

- We use the same 500 news to run our news similarity algorithm.
- The algorithm gives us cluster of news. Then we computer Confusion Matrix by comparing it with our Ground Truth.
- Now system generated clusters are compared with ground truth clusters and for each category Confusion Matrix is created. Confusion Matrix for category National is given below:

	3	1.	
\mathbf{P}	roc	110	ted
	TEC		UCU

		UNC	CN
Actual	UNC	TN = 91	FP = 1
	CN	FN = 3	TP = 5

Evaluation Measures

- F-measure
- Precision
- Recall

Evaluation Results

Category	Precision	Recall	F-Measure
International	0.62	0.5	0.55
National	0.83	0.63	0.71
Health	0.43	0.5	0.46
Business	0.26	0.35	0.3
Entertainment	0.29	0.33	0.31
Micro Average	0.45	0.46	0.46
Macro Average	0.48	0.46	0.47

Results

Sr#	Related News	Source
ande to to to-	بنگلادیش: بتاعت اسلامی کے ایک اور رہنا کے لیے سزائے موت	VOA
Cluster1	بنگلہ دلیق: ہاحتِ اسلامی کے رہنا اظہرالا سلام کو سزائے موت	BBC
	بنگلا دیش میں ٹر بوٹل نے جاعت اسلامی کے ایک اور رہنا کو سزائے موت سنادی	Express
	ذکی الرحمٰن لنحسوی کی نظربندی کا حکم معطل	BBC
Cluster2	ذکی الرحمٰن لنحموی کوایک مقدمے میں گرفتار کر لیا گیا	BBC
	فکی الرحمٰن کی فظربندی کا مکم معطل	VOA
	ذکی الرحمن لکھوی ایک اور مقدمے میں گرفتار	VOA
5	موٹر وے پولیس کی تخواہوں میں 20 فیصد اصافے کا اعلان	Jang
Cluster3	وزیراعظم کا موٹروے پولیں کی تخواہوں میں 20 فیصداضافے کا اعلان	Express
	وزیر اعظم کا موٹروے پولیس کی تخواہوں میں 20 فی صداحنانے کا اعلان	UrduPoint
Cluster4	4 ارب سال قبل مرتخ زمین بیساتها	BBC
	4ارب سال قبل ساره مربح کا ما تول زمین بیسا تھا، ناسا کا دعوی	NawaeWaqt
	4 ارب سال قبل مریخ کا ما تول زمین بیسا تھا	BBC
	سينٹ لوھيا: پاکستان کا بدف 243 رنز	BBC
Cluster5	سینٹ کوشیا: پاکسان کی پہلے ہیڈنگ	BBC
-	سينت لوهيا: ويسف الذي كابدف 230 رنز	BBC
	سينت لوشيا: پاکستان کابدف 262 رنز	BBC
	سينت ونسف: ويست انديز كالماكستان كوجيت كيلية 153 رز كابدت	Jang
	انسداد دہشت گردی کے قومی ادارے نیکٹا کو بحال کرنے کا فیصلہ	UrduPoint
Cluster6	انسداد دہشت گردی کے قومی ادارے " نیکٹا" کو بحال کرنے کا فیصلہ	Jang
	اانسداد دہشت گردی کے قومی ادارے کی فوری بحالی کی ہدایت	VOA
	انسداد دہشت گردی کے قومی ادارے 🛛 نیکٹا 🛛 کوبجال کرنے کا فیصلہ	GEO

Results (Cont.)

سائنس رصحت اشتارات	بۇنى	انثر فيتمنك	کھیل	انئر نيشنل	پاکستان	صفحة اول
اس جلیسی خبریں ۱۔ مشرق یوکرین لامانی میں •در جوں بلاک •			تفصيل	خبر کی		
2۔ یوکرین کا مشرقی حلاقوں میں کارردائیاں جاری رکھنے کا عزم 3۔ یوکرین بھران: مشرقی حلاقوں میں باغیوں کا ریفرینڈم		<i>پ</i> یثانی	،،عام لوگوں کی	ل ^ر ائی میں شدت	قی علاقوں میں	یوکرین کے مشر
4۔ یوکرین کے مشرقی علاقوں کوروس میں منم کرنے کا مطالبہ 5۔ روال ماہ مشرقی یوکرین میں لوائی سب سے ہلاکت خیزرہی:		-	an		/	
رپورٹ 6۔ یوکرین کا مشرقی علاقوں میں روسی فوج کے حکوں کا الزام		1				
7۔ مشرقی یوکرین کے نو حمر سپاہی 8۔ مشرقی علاقوں میں ریفرینڈم ایک 'ڈھونگ' : یوکرین 9۔ مشرقی علاقوں میں ریفرینڈم ایک 'نائک' : یوکرین	T					

Results (Cont.)

اشتارات	سائنس رصحت	يزنس	انثر فيتمنك	كهيل	انٹرنیشنل	پاکستان	صفحة اول
اس جیسی خبریں وگ کاران برقرار	1_پېچاب مېں سم			تفسيل	خبرک		
س افراد کو پیمانسی دے دی گئی پید 24 افراد میں ڈینگل دائر س کی تصدیق	2_معنجاب مين د			2n)	اد سموگ سے متا	12 سے زائدافر	•بنجاب م ی ں 00
بلگی وائرس سے متاثرہ افراد کی تعداد 54 ہو گئی زید69 افراد میں ڈینگی دائرس کی تصدیق				(jz			
اتخابات میں 50 افراد زخمی ب میں ملوث مزید پھافرادہ بنجاب سے گرفتار		کے اوقات میں	ر په نه دو سکې , صبح			خبارتازہ ترین ۔ 05	لاہور (اُردو پواتنٹ ا
ب حکومت یختلاف سردکول پر رسز کا آپریشن، مزید 250 افراد گرفتار		سموگ میں تخمی اور شام میں اس کی شدت بڑھ جاتی ہے , سموگ سے بارہ سوے زائد افراد متاثر ہو کر ہسپتال پہنچ گئے۔ دوسری طرف پنجاب کے مختلف علاقے دھند کی لپیٹ میں رہے , موٹر وے کئی مقامات پر بندر ہی۔ اس صورتحال					
مرچ آپریش؛200 افراد زیر حراست	10- پنجاب میں س	000			ے مرحد کی چیک پر یاط بر تنے کی اپیل ب		

Conclusion and Future Work

- We have an algorithm in place to integrate with newslink.pk to identify similar news automatically.
- Experimental evaluation shows an average micro average for Precision measure is 0.45 and the macro average for Precision is 0.48.
- In future, we intend to use LDA to identify similar news and compare it with our work.

THANKS!