

Model Input	Training Words			Training Vocab (Unique Words)		
	Count	% Full Text	Avg. per Article	Count	% Full Text	Avg. per Article
Full Text	5,032,755	100	684	87,541	100	264
Lead Paragraph	740,321	15	103	40,309	46	61
Headlines	57,752	1	8	11,150	13	6
Nouns	1,349,489	27	184	64,998	74	112
Lemmas	5,032,755	100	684	75,648	86	245

Table 1: Model Inputs

Desk	General Descriptor	Online Sections	Taxonomic Classifier
Foreign Desk	<ul style="list-style-type: none"> <li>• Immigration and Refugees</li> <li>• Jews</li> <li>• Music</li> <li>• Religion and Churches</li> </ul>	<ul style="list-style-type: none"> <li>• World</li> </ul>	<ul style="list-style-type: none"> <li>• Top/News/World/Europe</li> <li>• Top/News/World/Countries and Territories/Austria</li> <li>• Top/Features/Arts/Music</li> <li>• Etc. (8 others)</li> </ul>
Book Review Desk	<ul style="list-style-type: none"> <li>• Books and Literature</li> </ul>	<ul style="list-style-type: none"> <li>• Arts</li> <li>• Books</li> </ul>	<ul style="list-style-type: none"> <li>• Top/Feature/Arts</li> <li>• Top/Features/Books</li> <li>• Top/Features/Books/Book Reviews</li> </ul>
Classified		<ul style="list-style-type: none"> <li>• Paid Death Notices</li> </ul>	<ul style="list-style-type: none"> <li>• Top/Classifieds/Paid Death Notices</li> </ul>

Table 2: Possible Model Outputs

Cleaned Label Name	Original Label Name	Articles
book review	Book Review Desk	176
business & financial	Business World Magazine	1
	Business/Finance Desk	1
	Business/Financial Desk	628
	Business\Financial Desk	1
	E-Commerce	1
	Financial Desk	1106
	Financial Desk;	2
	Money and Business/Financial Desk	79
	SundayBusiness	14
cars	Automobiles	9
	Cars	4

Table 3: Label Cleaning Examples

Model Input	Model Type			Best Model
	MNB	LR	NN	
Full Text	<b>0.693</b>	0.742	0.	LR
Lead Paragraph	0.624	0.661	0.	LR
Headlines	0.502	0.545	0.	LR
Nouns	0.664	0.702	0.	LR
Lemmas	0.685	<b>0.743</b>	0.	LR
<b>Best Input</b>	Full Text	Lemmas	?	<b>Lemmas in LR</b>

Table 4: Model Accuracies on Test Data