







































































	Ex	ternal Me	asures	of Clu	Ister	Validity	y: Ent	tropy a	and Pu	urity	
	Table 5.9. K-means Clustering Results for LA Document Data Set										
(	Cluster	Entertainment	Financial	Foreign	Metro	National	Sports	Entropy	Purity		
	1	3	5	40	506	96	27	1.2270	0.7474		
	2	4	7	280	29	39	2	1.1472	0.7756		
	3	1	1	1	7	4	671	0.1813	0.9796		
	4	10	162	3	119	73	2	1.7487	0.4390		
	5	331	22	5	70	13	23	1.3976	0.7134		
	6	5	358	12	212	48	13	1.5523	0.5525		
	Total	354	555	341	943	273	738	1.1450	0.7203		
er	entropy For each cluster, the class distribution of the data is calculated first, i.e., for cluster $j$ we compute $p_{ij}$ , the 'probability' that a member of cluster $j$ belongs to class $i$ as follows: $p_{ij} = m_{ij}/m_j$ , where $m_j$ is the number of values in cluster $j$ and $m_{ij}$ is the number of values of class $i$ in cluster $j$ . Then using this class distribution, the entropy of each cluster $j$ is calculated using the standard formula $e_j = \sum_{i=1}^{L} p_{ij} \log_2 p_{ij}$ , where the $L$ is the number of classes. The total entropy for a set of clusters is calculated as the sum of the entropies of each cluster weighted by the size of each cluster, i.e., $e = \sum_{i=1}^{K} \frac{m}{m} e_j$ , where $m_j$ is the size of cluster j, K is the number of clusters, and $m$ is the total number of data points.										
թւ	<b>purity</b> Using the terminology derived for entropy, the purity of cluster $j$ , is given by $purity_j = \max p_{ij}$ and the overall purity of a clustering by $purity = \sum_{i=1}^{K} \frac{m_i}{m} purity_j$ .										
DЫ	From: Tan, Steinbach, Kumar, Introduction to Data Mining, McGraw Hill 2006										



