

	Doc1			Doc2			Doc3		
Step 0:  Raw data	Apple is a fruit which is red in colour.  Apple a day keeps doctor away.			Orange is a fruit which is orange in colour.  It is rich in vitamins.			Kiwi is a fruit which is green in colour.  It is a native fruit of new Zealand.		
Step 1:  Remove all stop words (like is, a, which, in, it, of, etc)	Apple fruit red colour apple day keeps doctor away			Orange fruit orange colour rich vitamins			Kiwi fruit green colour native fruit new zealand		
Step 2:  Lemmatize (i.e find the root word. Eg: root word of “seeing” is “see)	Keeps => keep  (no change in other words as they are root words already)  Apple fruit red colour apple day keep doctor away			Vitamins => vitamin  Orange fruit orange colour rich vitamin			(no change here)  Kiwi fruit green colour native fruit new zealand		
Step 3: Identify the frequency of each word within the document	Apple	2		Orange	2		kiwi	1	
	Fruit	1		fruit	1		fruit	2	
	Red	1		colour	1		green	1	
	Colour	1		rich	1		color	1	
	day	1		Vitamin	1		native	1	
	keep	1		total	6		new	1	
	doctor	1					zealand	1	
	away	1					total	8	
	total	9							
Step 4: Identify Term Frequency (TF) in each document  TF = (freq of word)/(no. of words in document)  Or  TF = 0.5 + 0.5*(freq of word)/(freq of most repeating word across documents)	term	TF of d1		TF of d2			TF of d3		
		0.222222							
	Apple	freq of word=2, No. of words in document=9, So TF=2/9=0.222222		0		0			
	Fruit	0.111111		0.166667		0.25			
	Red	0.111111		0		0			
	Colour	0.111111		0.166667		0.125			
	day	0.111111		0		0			
	keep	0.111111		0		0			
	doctor	0.111111		0		0			
	away	0.111111		0		0			
	orange	0		0.333333		0			
	rich	0		0.166667		0			
	vitamin	0		0.166667		0			
	kiwi	0		0		0.125			
	green	0		0		0.125			
	native	0		0		0.125			
	new	0		0		0.125			
	zealand	0		0		0.125			



<p>Step 7: <b>Find Cosine similarity between two documents</b></p> <p>Cosine Similarity = <math>\frac{(D1.D2)}{(\sqrt{\text{sum of squares of the TF IDF values of D1}} \cdot \sqrt{\text{sum of squares of the TF IDF values of D2}})}</math></p> <p>Where</p> <p><math>\sqrt{\text{sum of squares of the TF IDF values of D1}}</math></p> <p><math>\sqrt{\text{sum of squares of the TF IDF values of D2}}</math></p>	d1.d2	squares of d1	sqrt(squares sum of d1)	squares of d2	sqrt(squares sum of d1)	d1 . d2	Cosine Sim
	0	0.107747516	0.516839694	0	1.429534568	0.73884	0.0055698
	0.002057613	0.012345679		0.037037037			
	0	0.026936879		0			
	0.002057613	0.012345679		0.037037037			
	0	0.026936879		0			
	0	0.026936879		0			
	0	0.026936879		0			
	0	0.026936879		0			
	0	0		0.984747503			
	0	0		0.492373752			
	0	0		0.492373752			
	0	0		0			
	0	0		0			
	0	0		0			
	0	0		0			
<p>Similarly find Cosine Similarity for D1 and D3</p>	d1.d3	squares of d1	sqrt(squares sum of d1)	squares of d3	sqrt(squares sum of d3)	d1 . d3	Cosine Sim
	0	0.107747516	0.516839694	0	1.466955052	0.758181	0.0223895
	0.00308642	0.012345679		0.055555556			
	0	0.026936879		0			
	0.013888889	0.012345679		0.25			
	0	0.026936879		0			
	0	0.026936879		0			
	0	0.026936879		0			
	0	0.026936879		0			
	0	0		0			
	0	0		0			
	0	0		0			
	0	0		0.369280314			
	0	0		0.369280314			
	0	0		0.369280314			
	0	0		0.369280314			
<p>And find Cosine Similarity for D2 and D3</p>	d2.d3	squares of d2	sqrt(squares sum of d2)	squares of d3	sqrt(squares sum of d3)	d2 . d3	Cosine Sim

Step 8: **Identify most similar documents based on their cosine similarity values.**

Finally compare CosSim(d1, d2), CosSim(d1, d3) and CosSim(d2, d3).  
The 2 documents which have the highest CosSim value are the most similar documents.