

Counts

We've built a lot of fancy predictive models to, in part, determine which words (or parts of speech) are associated with a particular set of textual documents. What if we think about this problem directly, and just build a table showing how often a document occurs with a given label and a given term.

	Spam	Not Spam
!	178	40
no !	211	333

Counts

How strong of an association is there between the two? To see, let's compute the row and column counts.

	Spam	Not Spam	
!	178	40	218
no !	211	333	544
	389	373	762

Probabilities

Now, we will erase (for the moment) the counts and compute the probabilities associated with each category on its own.

	Spam	Not Spam		
!			218	28.6%
no !			544	71.4%
	389	373	762	
	51.0%	49.0%		

Expected Probabilities

Given just the totals, what the expected proportion of entries that should be in each cell? We can get these by multiplying the associated row and column probabilities.

	Spam	Not Spam		
!	0.510×0.286	0.490×0.286	218	28.6%
no !	0.510×0.714	0.490×0.714	544	71.4%
	389	373	762	
	51.0%	49.0%		

Expected Probabilities

Given just the totals, what the expected proportion of entries that should be in each cell? We can get these by multiplying the associated row and column probabilities.

	Spam	Not Spam		
!	14.5%	14.0%	218	28.6%
no !	36.4%	35.1%	544	71.4%
	389	373	762	
	51.0%	49.0%		

Expected Counts

Multiplying the probabilities by the number of documents (762) gives the expected counts.

	Spam	Not Spam		
!	111.2	106.8	218	28.6%
no !	277.5	266.6	544	71.4%
	389	373	762	
	51.0%	49.0%		

Measure Pe

Keep with this idea for a moment. With these proportions, we can compute the probability of observing the exact values (yes, it will be small) that we would get the observed data.

	Spam	Not Spam
!	14.5%	14.0%
no !	36.4%	35.1%

	Spam	Not Spam
!	178	40
no !	211	333

$$\text{Probability(Right | Left)} = P_e$$

Measuring Po

Similarly, we can compute the probability of observing the data given the observed proportions.

This will also be very small, but higher than the other number. The big question is: how much larger?

	Spam	Not Spam
!	23.4%	5.25%
no !	27.7%	43.7%

	Spam	Not Spam
!	178	40
no !	211	333

$$\text{Probability(Right | Left)} = P_o$$

G-Score

To measure the difference between these models, we compute what is called the g-score (or log-likelihood ratio). Higher values will correspond to words that are more strongly associated with a given label.

$$G = \log \left[\frac{\text{Probability}(\text{Left} \mid \text{Right}) = P_o}{\text{Probability}(\text{Left} \mid \text{Right}) = P_e} \right]$$

We can compute the G score for many terms and look at those that are the largest. We can extend this to multiclass classification by computing G scores for each specific category.