

Information Theory, Pattern Recognition and Neural Networks

HANDOUT 2 FEBRUARY 27, 2004

This is the assignment for week 8's supervision.

An old exam question

A spy would like you to write a computer program that recognises, given a small number of consecutive characters from the middle of a computer file, whether the file is an English-language document. Assuming that the two alternative hypotheses are that the file is an English-language document (\mathcal{H}_E), or that it is a random string of characters drawn from the same alphabet (\mathcal{H}_R), describe how you would solve this problem.

Estimate how many characters your method would need in order to work reasonably well.

Further questions

Maybe you would enjoy writing a program that implements your method?

How would your answers differ if instead the task were to distinguish

- (a) English from German?
- (b) English from Hsilgne (backwards English)?

[When I say German, let's assume unaccented German – German with all the umlauts stripped out.]

(Some facts about English and German are supplied below.)

LETTER FREQUENCIES OF ENGLISH AND GERMAN

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
English (e)	.07	.01	.03	.03	.10	.02	.02	.05	.06	.001	.006	.03	.02	.06
German (g)	.06	.02	.03	.04	.15	.01	.03	.04	.07	.002	.01	.03	.02	.08
	O	P	Q	R	S	T	U	V	W	X	Y	Z	-	
English (e)	.06	.02	.0009	.05	.05	.08	.02	.008	.02	.002	.01	.0008	.17	
German (g)	.02	.007	.0002	.06	.06	.05	.04	.006	.02	.0003	.0003	.01	.14	

The entropies of these two distributions are $H(\mathbf{e}) = 4.1$ bits; $H(\mathbf{g}) = 4.1$ bits; and the relative entropies between them are $D_{\text{KL}}(\mathbf{e}||\mathbf{g}) = 0.16$ bits and $D_{\text{KL}}(\mathbf{g}||\mathbf{e}) = 0.12$ bits. The relative entropies between the uniform distribution \mathbf{u} and the English distribution \mathbf{e} are $D_{\text{KL}}(\mathbf{e}||\mathbf{u}) \simeq 0.6$ bits and $D_{\text{KL}}(\mathbf{u}||\mathbf{e}) \simeq 1$ bits.

How well calibrated are your estimates of uncertainty?

According to BBC News, the Earth was almost put on impact alert by some astronomers who, on 13 January 2004, reckoned a 30m object, later designated 2004 AS1, had a one-in-four chance of hitting the planet within 36 hours. 2004 AS1 in fact passed the Earth at a distance of about 12 million km – 32 times the Earth-Moon distance, posing no danger to us whatsoever. Clearly the correct assessment of the probability of earth being hit should have been more like one in 10^7 .

This sheet is intended to give insight into the accuracy of our own probability estimates.

Give a 94% confidence interval for the following quantities. Give the tightest* interval you can, while remaining 94% sure that the true value is in the interval.

	Quantity	Lower bound	Guess	Upper bound	Ratio	Score
1	Mass of textbook (g)					
2	Population of Britain					
3	Population of Turkey					
4	Population of Luxembourg					
5	Number of British MEPs					
6	Starting pay of University Lecturer					
7	Parliamentary salary of MP					
8	Council tax, South Cambs. (£/house)					
9	Fraction of central government expenditure that goes to ‘Defence’					
10	UK prison population (as fraction of whole)					
11	Number of USA nuclear warheads					
12	Distance to sun (miles)					
13	Mean radius of earth (km)					
14	Avogadro’s number (one mole)					
15	Speed of light (m s^{-1})					
16	Standard atmospheric pressure (pa)					

* My definition of ‘tightest’ is that an interval (l, u) gets tighter if the ratio u/l is reduced. My measure for how tight an interval is is $\log_2(\log_2 u/l)$ – the smaller, the better.

This score is a measure of how useful a compression service your confidence intervals provide. If Fred is about to learn the value of x and needs to compress it, and doesn’t know x to within a factor of 10^5 , then if you supply a confidence interval of width 1% (thus replacing 10^5 by 1.01) then Fred can save 10 bits by trusting your 94% interval; on the other hand, if the true value of x falls outside your interval – which you expect to happen 6% of the time – then Fred will pay a *penalty* of $\log_2 1/0.06 = 4$ bits.

In summary, you can score your answers by adding +4 bits every time the true value falls outside your range; and a score of $\log_2(\log_2 \text{Ratio}) - 4$ bits if it’s inside your range.

The resulting total is hopefully a negative number, representing the bit-saving that Fred gets by using your knowledge to compress $\{x\}$ instead of using his knowledge (a range of size 10^5 for every quantity).

Ratio	$\log_2(\log_2 \text{Ratio})$	$\log_2(\log_2 \text{Ratio}) - 4$
1.01	-6	-10
1.02	-5	-9
1.04	-4	-8
1.09	-3	-7
1.19	-2	-6
1.41	-1	-5
2	0	-4
4	1	-3
16	2	-2
250	3	-1
100,000	4	0