

Synonymy and corpus work: On *almost* and *nearly*

Göran Kjellmer
Göteborg University

Dictionaries are indispensable tools for a language learner. They will tell him what words mean and how they are distinguished semantically. What the dictionaries do not always do is specify how words are used idiomatically by native speakers, and in omitting to do so they can sometimes be misleading in that the information they offer may be necessary without being sufficient. Consider the case of *almost* and *nearly*.

Almost and *nearly* are close synonyms, as is clear from a few dictionary definitions:

almost all but; very nearly

nearly almost

(COD 1990)

almost very nearly but not completely

nearly almost, but not quite or not completely

(LDOCE 1995)

almost nearly but not quite

nearly almost or not completely

(CIDE 1995)

almost not quite; very nearly

nearly very close to; almost

(NODE 1998)

You use **almost** to indicate that something is not completely the case but is nearly the case. **Nearly** is used to indicate that something is not quite the case, or not completely the case.

(Cobuild 2001)

It appears from those definitions that the words are so synonymous that they are sometimes defined in terms of each other. If they are so close in meaning, one may well wonder if there are any differences between them in the way they are used. I will here discuss three aspects, viz:

- their frequency
- their style and text type preference
- their collocability

In order to do so I will make use of the CobuildDirect Corpus, where such aspects are convenient to deal with.

A. Frequency. *Almost* is more than twice (2.33 times) as frequent as *nearly* in CobuildDirect; it has 15,536 occurrences, while *nearly* has 6,666. Whether this difference in frequency between the words is significant or not (in a non-statistical sense) when it comes to their use is too early to say. But the difference could suggest that *nearly* is a more select and specialised word than *almost*.

B. Text type preference. CobuildDirect consists of twelve subcorpora, taken from different text types. The distribution of *almost* and *nearly* over those text types might shed some light on their characteristics of usage.

Table 1: Text type preference of *almost* and *nearly*

Almost:

usbooks	2256	401.0/million
ukbooks	1939	362.1/million
times	2027	351.7/million
npr	1039	332.0/million
oznews	1538	288.1/million
ukmags	1337	272.7/million
today	1422	270.9/million
bbc	669	256.3/million
sunnw	1374	235.9/million
ukephem	536	171.6/million
usephem	177	144.5/million
ukspok	1222	131.8/million

Nearly:

npr	669	213.8/million
bbc	532	203.8/million
today	814	155.1/million

usbooks	794	141.1/million
times	758	131.5/million
sunnow	666	114.3/million
usephem	122	99.6/million
oznews	489	91.6/million
ukmags	449	91.6/million
ukbooks	461	86.1/million
ukspok	690	74.4/million
ukephem	222	71.1/million

The abbreviations for the sources are to be interpreted thus:

npr	= US National Public Radio broadcasts
today	= UK Today newspaper
times	= UK Times newspaper
usbooks	= US books; fiction & non-fiction
oznews	= Australian newspapers
bbc	= BBC World Service radio broadcasts
usephem	= US ephemera (leaflets, adverts, etc.)
ukmags	= UK magazines
sunnow	= UK Sun newspaper
ukspok	= UK transcribed informal speech
ukbooks	= UK books; fiction & non-fiction
ukephem	= UK ephemera (leaflets, adverts, etc.)

The middle column gives the raw frequencies of the words in the subcorpora, and the rightmost column converts those figures to comparable frequencies per one million words.

The table shows that *almost* prefers literary styles of writing (usbooks, ukbooks, times) and avoids more popular text types (sunnow, ukephem, usephem, ukspok), whereas *nearly* is more strongly favoured in the news media (npr, bbc, today). Neither of them is used much in spoken British English (ukspok).

C. Collocability. The corpus provides a facility by means of which the most significant collocates of any word in the corpus can be discovered. Four words on either side of a word are seen as its collocates. T-score calculations will then indicate their significance and sort the words accordingly. The most significant collocates of *almost* and *nearly* are given in Table 2.

Table 2: Collocates of *almost* and *nearly*

Collocates of <i>almost</i>			Collocates of <i>nearly</i>		
	n	T-score		n	T-score
<i>certainly</i>	650	24.637057	<i>years</i>	681	23.963591
<i>as</i>	1569	23.223282	<i>000</i>	354	17.239044
<i>every</i>	486	19.473519	<i>all</i>	533	16.802179
<i>it</i>	2164	19.356305	<i>million</i>	287	15.773051
<i>years</i>	562	18.235811	<i>half</i>	285	15.633906
<i>impossible</i>	269	16.002639	<i>two</i>	391	15.562322
<i>was</i>	1558	15.886701	<i>ago</i>	217	13.558705
<i>an</i>	800	15.398975	<i>three</i>	252	12.840255
<i>all</i>	707	13.871836	<i>every</i>	201	12.461522
<i>immediately</i>	208	13.724134	<i>hundred</i>	163	12.138026
<i>certain</i>	224	13.649910	<i>pound</i>	235	11.784722
<i>half</i>	275	13.621994	<i>thousand</i>	141	11.349662
<i>000</i>	305	13.507498	<i>dollar</i>	173	11.160020
<i>entirely</i>	149	11.767122	<i>cent</i>	157	11.070751
<i>million</i>	212	11.392851	<i>as</i>	534	11.056001
<i>cent</i>	199	11.085978	<i>per</i>	160	10.858025
<i>ago</i>	195	11.082241	<i>after</i>	240	10.485561
<i>is</i>	1509	10.987927	<i>four</i>	162	10.352990
<i>always</i>	211	10.592463	<i>year</i>	219	10.028319
<i>dollar</i>	210	10.275591	<i>for</i>	708	9.756185
<i>any</i>	303	10.048554	<i>months</i>	126	9.731231
<i>exclusively</i>	104	10.010189	<i>always</i>	128	9.146864
<i>certainly</i>	650	24.637057	<i>five</i>	131	9.061732
<i>as</i>	1569	23.223282	<i>six</i>	112	8.727735
<i>every</i>	486	19.473519	<i>30</i>	100	8.583709
<i>it</i>	2164	19.356305	<i>20</i>	100	8.492225
<i>years</i>	562	18.235811	<i>points</i>	86	8.350813
<i>impossible</i>	269	16.002639	<i>percent</i>	82	8.308806
<i>was</i>	1558	15.886701	<i>hours</i>	90	8.261594
<i>an</i>	800	15.398975	<i>40</i>	81	8.235720
<i>all</i>	707	13.871836	<i>billion</i>	75	8.178398
<i>years</i>	562	18.235811	<i>200</i>	71	7.956700
<i>impossible</i>	269	16.002639	<i>1</i>	136	7.912442
<i>was</i>	1558	15.886701	<i>died</i>	70	7.407585

<i>an</i>	800	15.398975	<i>50</i>	77	7.369675
<i>all</i>	707	13.871836	<i>300</i>	59	7.320331
<i>immediately</i>	208	13.724134	<i>2</i>	111	7.254769
<i>certain</i>	224	13.649910	<i>twenty</i>	67	7.203400
<i>half</i>	275	13.621994	<i>lost</i>	72	7.113234
<i>000</i>	305	13.507498	<i>quarter</i>	56	6.954470
<i>entirely</i>	149	11.767122	<i>weeks</i>	69	6.920944
<i>million</i>	212	11.392851	<i>killed</i>	59	6.845009
<i>cent</i>	199	11.085978			
<i>ago</i>	195	11.082241			
<i>is</i>	1509	10.987927			
<i>always</i>	211	10.592463			
<i>dollar</i>	210	10.275591			
<i>any</i>	303	10.048554			
<i>exclusively</i>	104	10.010189			
<i>like</i>	421	9.506938			
<i>per</i>	171	9.038828			
<i>everything</i>	127	8.795240			
<i>seems</i>	117	8.679364			
<i>seemed</i>	111	8.668655			
<i>two</i>	351	8.375584			
<i>everyone</i>	104	8.321448			
<i>had</i>	614	8.220298			
<i>anything</i>	134	8.190833			
<i>percent</i>	95	8.130223			
<i>identical</i>	67	7.996534			
<i>completely</i>	80	7.844077			
<i>twice</i>	78	7.771227			
<i>could</i>	288	7.512932			
<i>daily</i>	73	7.342104			
<i>40</i>	81	7.218745			
<i>three</i>	218	7.161602			
<i>after</i>	306	7.159486			
<i>she</i>	492	7.119327			
<i>has</i>	529	7.095944			

It appears from the table that typical collocates of *almost* and *nearly* are adverbs, adjectives, pronouns, prepositions, nouns and numerals in varying proportions. If we focus on the postcontexts of our words, Table 3 gives the distribution of the parts of speech immediately following them.

Table 3: Part-of-Speech-sorted words immediately following *almost* and *nearly*

	Adverbs	Adjectives	Pronouns	Verbs	Nouns	Numerals	Prep
<i>Almost</i>	2921	2743	561	2443	689/392	1840/2144	731
<i>Nearly</i>	395	366	106	1012	486/150	2633/2969	91
<i>Almost</i>	7.39	7.49	5.29	2.41	1.42/2.61	0.70/0.72	8.03
<i>Nearly</i>							

Almost is 2.33 times as frequent as *nearly*. If *almost* is much more, or much less frequent than 2.33 times the frequency of *nearly* in a given category, this would consequently be interesting.

Some preliminary conclusions are these. Adverbs, adjectives, pronouns and prepositions are typical postoccurring collocates of *almost*; nouns and numerals are typical postoccurring collocates of *nearly*. (The collocates preceding *almost* and *nearly* are most of the time, or at least very frequently, syntactically unrelated.)

However, some of these figures are misleading. Sums of money like \$100 or £100 are given as “$ 100” and “£ 100” and classified as nouns. In this context it might be reasonable to regard them as numerals, and if they are treated as such, and if the difference in overall occurrence is taken into account, the difference between *almost* and *nearly* with regard to a following noun becomes uninteresting. (The adjusted figures are those following the slash.)

That the PoS categories are too general and conceal some important facts can be seen from the following experiment. There are 7.39 times as many adverbs following *almost* as there are following *nearly*. But what are those adverbs? Table 4 gives the beginnings of the lists “*almost* + adverb” and “*nearly* + adverb”:

Table 4: *Almost/nearly* + adverb (RB)

a result-oriented judge who sided	almost always with the government against
matches twice a day, he would play	almost no more than an hour and a half and
if the US makes the final and plays	almost undoubtedly Australians in Saint
Algeria and Jordan met quietly,	almost obscurely in Rabat. All the

the Life of Andy Warhol,' set has laid out some scenarios which did not foresee that this would lead phone call which was followed by the PLO that the US would names, and one of those teams will markets. These people must rely he fell very fast. And people were better you get, you know, and I'm

The two countries went to war s happened. It's--it has happened Lyman: Wirtz's mythical town is the road. One w--will be crushed, in with firm instructions and they go in with firm instructions and broadcast. When they are they're a realist, but audiences and critics Kuwait. Such a resolution would non-entities, whose names were

in the wry music of Algeria and the US Public Interest Research Group, distribution rights, though not year, when Stemple earned makes furniture that could be put hospitals. He found men were black and white, and--and it's not immediate reaction would not be combined German team won't be I can.' By now her eyes burned fail ridiculously. But I don't feel that six-year-old children were as an insider or an outsider is not moves, which they say do not go to the ancient Egyptians, he wasn't

Editor The parties are not US government is not spending student athletes graduate, a rate says the prosecutor did not go the kidnap-for-ransom gangs

almost entirely to the music of the early almost certainly would draw US military almost immediately to a severe meat almost immediately by the bank calling to almost certainly have to veto. Yesterday almost surely be the national champion almost entirely on state food stores for almost deliberately staying away from almost as old as Sonny Boy Williamson was almost immediately, and when they finally almost everywhere. There's no parallel to almost as densely populated as Lake almost unidentifiably, and the other will almost intentionally, it seems to me almost intentionally, it seems to me almost always used as propaganda. almost unanimously saw him as a romantic almost certainly be defeated, sources say. almost as long as the parts they play. In

nearly as famous as Oum Kalsoum but very nearly twice as many banks are charging nearly as many as they had hoped for. nearly twice that amount. Other GM nearly anywhere. That's fortunate, since nearly twice as likely to undergo nearly as polished as this <p> Simon: Mm-nearly as sympathetic and moreover that a nearly as good as the East German team was nearly too high. I could see she was nearly as foolish as the time I bought the nearly as familiar with the Old Joe nearly as important as proving to voters nearly far enough, will be enough for nearly as important as an earlier king, nearly as powerful as they once were, but nearly enough on industrial research and nearly twice the NCAA average. But nearly far enough. She says all seven of nearly always got what they asked for,

2nd OPERAGOER: Oh, nearly, to a strict budget which is not
HAZEL: Christmas is nearly, yes laughs <p> FRANKL: May nearly as big as it was last year <p> At nearly here and I'm exhausted already.

As we can see, *almost* is typically followed by manner adverbs (*obscurely, intentionally*), time adverbs (*always, immediately*) and sentence adverbs (*undoubtedly, certainly*), whereas *nearly* typically occurs in the construction *not + nearly + as*, where *as* is classified as an adverb. There is thus a qualitative as well as a quantitative difference between *almost* and *nearly* in this respect.

Summary and conclusion. The differences found between *almost* and *nearly* are these:

- *almost* is much more frequent than *nearly* and is therefore likely to be less specialised than the latter;
- *almost* occurs more in literary styles of writing than in popular text types, whereas *nearly* is a preferred word in the news media. Neither of them is used much in the spoken language;
- their collocations distinguish the two words sharply, so that *almost* is characteristically followed by adverbs (*almost certainly*), adjectives (*almost impossible*), pronouns (*almost anything*) and prepositions (*almost by definition*), and *nearly* is equally characteristically followed by numerals (*nearly 200 people*).

These three aspects can be seen to be interrelated. *Nearly*, which occurs more often in the news media, where precision and factual information are more focused than in literary styles, is to some degree specialised in that it is preferably used to modify precise figures.

So far, from being the next-to-interchangeable synonyms that dictionaries could lead us to imagine, *almost* and *nearly* turn out, on closer inspection, to be partly overlapping but in important respects clearly contrasting words. Even if they are closely related in meaning, corpus studies can show that they are used differently by idiomatic speakers of English. Identity or near-identity of dictionary definitions does not guarantee identity or near-identity of usage.

References

CIDE = Procter, Paul (ed.). 1995. *Cambridge International Dictionary of English*. Cambridge University Press.

- Cobuild = Sinclair, John (editor-in-chief) 2001. *Collins COBUILD English Dictionary for Advanced Learners*. 3rd ed. Glasgow: HarperCollins.
- COD = Allen, R.E. (ed.). 1990. *The Concise Oxford Dictionary of the English Language*. 8th ed. Oxford: Clarendon.
- LDOCE = Summers, Della (ed.). 1995. *Longman Dictionary of Contemporary English*. 3rd ed. Harlow, Essex: Longman.
- NODE = Pearsall, Judy (ed.). 1998. *The New Oxford Dictionary of English*. Oxford: Clarendon.