

Combining Formal and Usage based Theories with Data Science Techniques in Measuring the Development of Syntactic Complexity in Written Production

https://github.com/ELI-Data-Mining-Group/

Alan Juffs and Na-Rae Han juffs@pitt.edu

Department of Linguistics

Kenneth P. Dietrich School of Arts and Sciences

Collaborators

- Ben Naismith presenting at TESOL
- Daniel Zheng

- Funding
- Pittsburgh Science of Learning Center
- http://www.learnlab.org



Principled Eclecticism in SLA research

- 1. Data Science slice and dice large data sets
- 2. Formal linguistics: syntactic categories
- 3. Corpus linguistics: frequency counts
- 4. Psycholinguistics: experiments in processing
- Each has a contribution to make: not mutually exclusive and can be combined
- 'Closing the loop' with educators



1. Data Science Tools

- http://www.pitt.edu/~naraehan/ling1340/Lecture1.pdf
- Python: can import very large .csv files
 - Pandas Data frames to organize and label data
 - Numpy mathematical operations/statistics
 - Competing with R
 - NLTK lexical tools (well known)
 - Matplotlib data visualization tools
- ELI tool kit developed in Python Daniel Zheng
 - Lemmatizer Someya list (!)
 - D measure of diversity
 - Advanced Guiraud sophistication growth of less frequent lexical items dependent on L1/culture
 - (Juffs, 2019; Naismith et al. 2018)



Import Data: tokenize, lemmatize, POS tag

In [2]: combo_dir = "../../Data-Archive/"

```
# Read two pickles as separate DFs
combol = pd.read_pickle(combo_dir + 'combo_df_lA_pt1.pkl')
```

```
combo2 = pd.read_pickle(combo_dir + 'combo_df_1A_pt2.pkl')
```

```
# Unpickle two DFs, and then concatenate along axis 1 (y axis)
combo_df = pd.concat([combo1, combo2], axis=1)
combo_df.head()
```

Voila! One big combo_df, ready to dice.

Out[2]:

question_id anon_id user_file_id text class_code level_id native_language version toks_re toks_re_len toks_nItk toks_pos lemmas

'er_id													
1	5	eq0	7505	I met my friend Nife while I was studying in a	g	4	Arabic	1	[i, met, my, friend, nife, while, i, was, stud	177	[I, met, my, friend, Nife, while, I, was, stud	[(I, PRP), (met, VBD), (my, PRP\$), (friend, NN	[i, meet, my, friend, nife, while, i, be, stud
2	5	am8	7506	Ten years ago, I met a women on the train betw	g	4	Thai	1	[ten, years, ago, i, met, a, women, on, the, t	137	[Ten, years, ago, ,, l, met, a, women, on, the	[(Ten, CD), (years, NNS), (ago, RB), (,, ,), ([ten, year, ago, ,, i, meet, a, woman, on, the
3	12	dk5	7507	In my country we usually don't use tea bags. F	w	4	Turkish	1	[in, my, country, we, usually, don, t, use, te	63	[In, my, country, we, usually, do, n't, use, t	[(In, IN), (my, PRP\$), (country, NN), (we, PRP	[in, my, country, we, usually, do, n't, use, t
4	13	dk5	7507	l organized the instructions by time.	w	4	Turkish	1	[i, organized, the, instructions, by, time]	6	[l, organized, the, instructions, by, time, .]	[(I, PRP), (organized, VBD), (the, DT), (instr	[i, organize, the, instruction, by, time, .]
5	12	ad1	7508	First, prepare a port, loose tea, and cup.\nSe	w AAAL 2019	4 Marcl	Korean	1	[first, prepare, a, port, loose, tea, and, cup	59	[First, ,, prepare, a, port, ,, loose, tea, ,,	[(First, RB), (,, ,), (prepare, VB), (a, DT),	[first, ,, prepare, a, port, ,, loose, tea, 5

2. Formal Linguistics and Corpus Analysis: Making Predictions: what to look for

- White (1987)
 - Dative alternation: not restricted to input (c.f. C. L. Baker (1979)
 - 'Mummy, open Hadwen the door'
- Zobl (1989)
 - Unaccusative/ Unergative verbs: passive overused with unaccusatives
- Schwartz & Sprouse (1996) L2 German (Cevdet)
- Lardière's work
 - Feature reassembly (Patti) and L2 Korean plurals (e.g., Hwang & Lardiere, 2013)



Lexical and Functional Categories

- Verbs Pinker (1989); Juffs (1996)
- Tense and morphology:
 - Prevost & White (2000): French L2 finite/nonfinite
- if, whether CP

Sam wonders whether it will rain.

- the, that, 's DP
 - Sam's book
 - *The Sam's book



3. Corpora and Usage-Based Approaches

- Frequency ranks- work in NLP .. Crossley, Kyle, Jarvis – Sunday's colloquium
- N. Ellis (2016, pp. 44-46)
 - Ortega (2001): Longitudinal data needed
 - Few longitudinal L2 corpora available
- Debates about which measures of association are most relevant to acquisition and processing
 - Must rely on 'hundreds of millions of words to approximate usage' (Ellis, 2016, p. 44).
- Too pessimistic? Predict (some) development based on formal theories?



Bley-Vroman (2002)

- In the meaning-based approach, the statistical structure of the language can affect the development of linguistic knowledge (for example, by influencing acquisition order or providing evidence for developing grammars);
- However, *linguistic* knowledge is NOT itself knowledge of the statistical structure of language
- We need to consider what learners want to communicate as well
- Argument elaborated on by Yang (2008)





4. Psycholinguistics research: norming: Kennison (1999)

VSam knows [NP the answer]. VSam knows [CP the answer is correct] ?Sam supposed [NP the answer].VSam supposed [CP the answer is correct]

Verb	Difficulty for learners? COCA Frequency?	NP %	CP %
consider		93	0
suggest		32	59
explain		82	4
realize			СР
admit		14	42
deny		78	11
conclude		25	63
recommend		55	44
suppose			СР





- PELIC Large and longitudinal. "In the wild!"
- <u>https://github.com/ELI-Data-Mining-Group/Pitt-ELI-Corpus</u>



Number of texts

Number of Texts > 10 Words by L1 and Level



Pitt IEP Levels and Cut Scores

(No beginners)

Level	ept Combined	ELI Writing	CEFR Equiv.
2. High Beginner	28-37	2.1-2.9	A1 (Breakthrough)
3. Low Intermediate	38 - 47	3.1 – 3.9	A2-B1 (Waystage)
4. Intermediate	48 - 59	4.0 - 4.9	B1 (Threshold)
5. High Intermediate	60 - 68	5.0 – 5.9	B2 - edge of C1 (Vantage)
6. Low Advanced	69 +	6.0	Low C1 (Effective)
		"Writing used in borderline cases."	



All Written Data: 'first version'

Level	Token Count	Per Mil Multipl
3	524,137	1.9
4	1,628,232	0.61
5	1,462,346	0.68
Total	3,614,715	



Data: Re-Token Counts x L1 and Level

Level	Arabic	Chinese	Korean	Total
3	157,727	71,431	74,813	303,971
4	398,333	278,964	296,291	973,588
5	364,614	298,430	259,615	922,659
Total	920,674	648,825	630,719	2,200,218



Create Frequency Ranks: L1 x Level

COCA Word Frequency Ranks

5	Rank	Word	Part of speec	Frequency
5				_
7	1	the	а	22038615
3	2	be	V	12545825
)	3	and	С	10741073
)	4	of	i	10343885
L	5	а	а	10144200
2	6	in	i	6996437
3	7	to	t	6332195
ļ	8	have	V	4303955
5	9	to	i	3856916
5	10	it	р	3872477
7	11	I	р	3978265
3	12	that	С	3430996
)	13	for	i	3281454
)	14	you	р	3081151
L	15	he	р	2909254
2	16	with	i	2683014
3	17	on	i	2485306
1	18	do	V	2573587
5	19	say	V	1915138
5	20	this	d	1885366
7	21	they	р	1865580
3	22	at	i	1767638
)	23	but	C	1776767
				1000005

Korean Level 3 – freq. per mil

Ζ	А	В	С
1		lemma	count
2	1	be	44083
3	2	i	35716
4	3	the	33163
5	4	to	30623
6	5	а	29139
7	6	and	22737
8	7	of	18847
9	8	in	17337
LO	9	have	17056
Ι1	10	you	13674
۱2	11	for	11054
L3	12	my	10600
L4	13	that	9517
٤5	14	it	8769
16	15	do	7860
L7	16	he	7525
18	17	can	6550
L9	18	n't	6376
20	19	we	5988
21	20	they	5788





Research Questions

- To what extent do *selected* lexical items in the COCA frequency ranking reflect frequency rankings of those words in ESL students written output?
- What are the frequency ranks of single Functional Category words (associated with morpho-syntactic complexity) in the most frequent 3000 lemmas of L2 written output?
- What are the frequency rankings of verbs requiring complex syntax? E.g., 'know' and 'suppose'
- Can these results inform classroom practice = closing the loop?



Functional Categories: All Data



Development of 'whether'

Est. Frequency per Million of 'whether'



Verbs: Kennison (1999)

	COCA Rank	NP	СР	BNC-COCA-25	IEP List?
consider	395	Х		1000	considerable
suggest	431		Х	1000	X
explain	481	Х		1000	X
realize	621		X	1000	X
admit	1093		Х	1000	X
deny	1413	Х		2000	V
conclude	1680		Х	3000	\checkmark
recommend	1699	Х		2000	X
suppose	2118		X	1000	X



Verbs: estimated frequency per million for comparison



















Summary

Changes over Level in Verbs favoring NP and CP: est. Frequency per Million Words



Discussion

- Many words in ESL production are predictable based on frequency in COCA – confirms usefulness of frequency bands
- Selecting lexical items based on theoretical part of speech status shows that frequency is not the only determinant of use in written output
- Syntactic complexity: markers of complex T-units increase with level: verbs requiring complex semantics --> CP selection (Grimshaw, 1981)
 - Topic choice/what learners want to say: 'concluding' vs. 'suggesting'
- Words chosen on formal theoretical/experimental psycholinguistic grounds are a promising direction a proxy for measuring syntactic development: now need to check whether verbs are used with NP or CP in reality



Conclusion

- Data Science tools +theory permits principled exploration of 'big' learner data sets
- Add the use of formal theories, not just usage based ones
- Include insights from psycholinguistics
- And finally



Close the loop!

- 1. Beginning of the loop students and teachers provide data
- 2. Researchers analyze the data
- 3. Researchers close the loop
- → *Discuss* with IEP teachers, e.g.,
 - which words to focus on in valuable class time (when the learners need to acquire 8-9000 words).
 - Less 'considering'; more 'suggesting', 'admitting', 'denying'
 - Some learners may need more focus on key clause types, e.g., Arabic speakers: 'whether';
 - All learners could use focus on CP verbs



Robert Ochsner

Ochsner, R. (1979). A poetics of second language acquisition. *Language Learning, 29*(1),

53-80.

Man is a complex symbol-user, and language is the basic tool he uses to symbolize everything. Language *is* a very big topic. SLA comprises obviously as big a field of study. It should not surprise us then to find at least two possible ways to investigate a first or second language. If we study it as a simple entity, a "body of facts," we can develop by controlled experimentation a good sense of the *thing* that language is. But if we only look at its physical use, we miss understanding the talker's symbols, his/her metaphorical use of language. It is possible to study either the language *thing* or the language *metaphor;* our research attitude determines this view or way of looking. Therefore, at another research level, we must look at our looking, and then through our point-of-view. And then *at* and *through*, perspective and reality in constant oscillation (Lanham 1977). For the researcher shapes his research findings, a fact that has been experimentally, and ironically, shown (Rosenthal 1977).

This at/through perspective is itself, remember, a metaphor—a kind of applied poetics. Paula Kurman (1977), a speech therapist, gives the nicest statement of what I mean, in practice, by this "bilingual" attitude. In sum, I can only steal her remarks:

It is useful to read of developments in fields other than one's own.

It is useful to "borrow" question structures and research techniques from other fields.

Is is useful to seek out perspectives of those whose work interests are different from one's own.

It is useful *not* to restrict such solicitations to the scientific disciplines. Artists, business people, and alert children have interesting perceptual frameworks, too. Remember that it was a child, as yet not thoroughly enmeshed in the social construct of reality of his elders, who saw and called the Emperor naked.







- Compose meaningful sentences and paragraphs that focus on a central idea with appropriate support and conclusion
- Introduce the concept that writing is a process
- Express ideas in writing to the reader in as clear a way as possible
- Increase fluency in writing



- BNC 3-4000 level words
- 'Students will produce medium-length, original written texts (≤ 500 words) responding to information on personal, practical, social, and general academic topics.'
- Level 4 topics: 'Process' (common examples were 'recipes' and 'how to find an apartment');
 'Classification' (e.g., types of doctors, festivals, jobs, lists of reasons). 'Cause-effect': one frequent topic was of 'causes of happiness'.



- BNC 5000 list + Coxhead core
- 'Students will produce medium-length and long, original written texts (500-2,000 words) on personal, practical, social, and general academic topics.'
- Level 5:
 - 'Explanation'
 - (e.g., 'how to learn English', 'how to stay healthy', 'the effects of a bad diet');
 - 'Narratives';
 - 'Argument/persuasive essay' that presents a point of view and supports it (e.g., 'euthanasia', 'the death penalty', 'pollution'; 'same-sex marriage');
 - 'Comparison/contrast essay'
 - (e.g., 'town or city living', 'your home town vs. Pittsburgh', 'Macintosh vs. PC computers').
 - 'Example Essay'
 - illustrate a case, e.g., 'education in the ELI', etc.





