A Readable Read: Automatic Assessment of Language Learning Materials based on Linguistic Complexity (2016) by Ildiko Pilan, Sowmya Vajjala, Elena Volodina

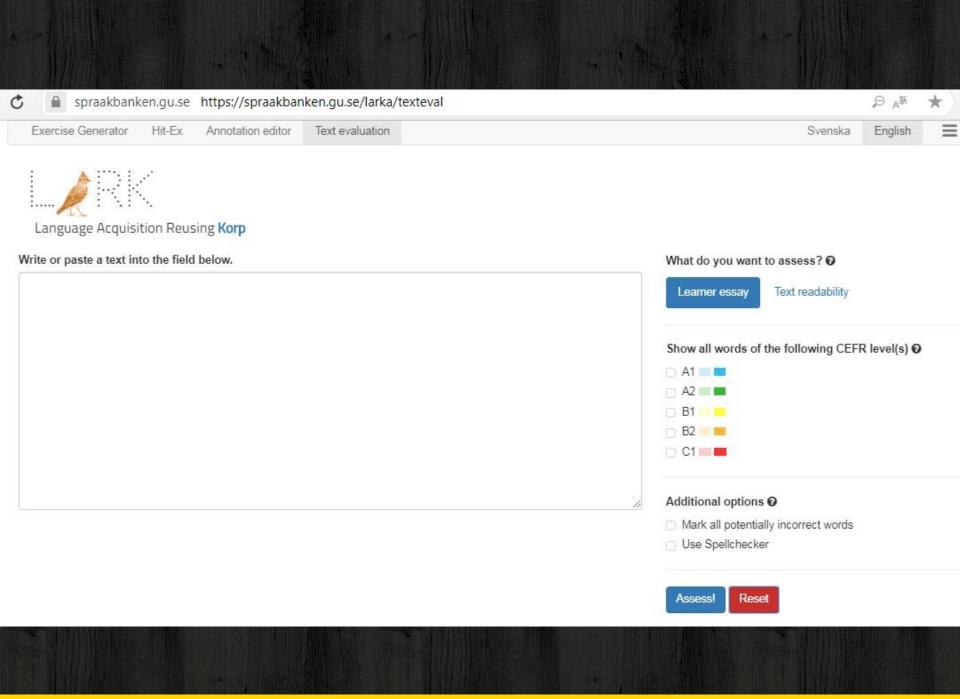
Цель: предсказать сложность текста на шведском для изучения его иностранцами

Результат: модель с точностью 81.3% и F-мерой 0.8

Läsbarthetsindex, 'Readability index'

$$LIX = \frac{A}{B} + \frac{C \times 100}{A}$$

- А количество слов в тексте,
- В количество предложений в тексте,
- С количество слов длиннее 6 букв
- < 30 детские тексты, 30–40 газетные статьи, 40–50 – журнальные статьи, 50–60 – официальные тесты, > 60 – законы



Данные

- COCTAILL corpus
- Уровни А1-С1
- 12 книг (4 издательства)
- Тексты аннотированы: POS и dependency tags
- 867 текстов
- Основаны на диалогах

Данные

CEFR	Books	Publ.	Texts	Mean num. sentences
A 1	4	3	49	14.0
A2	4	3	157	13.8
B1	5	3	258	17.9
B2	4	3	288	26.6
C 1	2	2	115	42.1
Total	12	4	867	

	r. Feature	Name Nr.	Feature Name
	Length-b	ised	Morphological
15/19/19/19/19/19/19/19/19/19/19/19/19/19/	Sentence	length 30	Modal verbs to verbs
	Average to	ten length 31	Particle INCSC
	Extra-lon	g words 32	3SG pronoun INCSC
	Number of	characters 33	Punctuation INCSC
	LI	X 34	Subjunction INCSC
	Lexico	I 35	S-verb INCSC
	A1 lemm	a INCSC 36	S-verbs to verbs
国际和国际企业的 国际包含国际国际的	A2 lemm	a IncSc 37	Adjective INCSC
	B1 lemm	a IncSc 38	Adjective variation
	B2 lemm	INCSC 39	Adverb INCSC
	C1 lemm	a INCSC 40	Adverb variation
	C2 lemm	INCSC 41	Noun INCSC
THE PROPERTY OF THE PROPERTY O	2 Difficult we	ord INCSC 42	Noun variation
	B Difficult noun a	nd verb INCSC 43	Verb IncSc
	Out-of-Ke	ly IncSc 44	Verb variation
	Missing lemma	form INCSC 45	Nominal ratio
		100 CA 100 C	Nouns to verbs
	Syntac		Function word INCSC
Features	7 Average deper	dency length 48	Lexical words to non-lexical words
	B Dependency arc	s longer than 5 49	Lexical words to all tokens
	Longest dependence	y from root node 50	Neuter gender noun INCSC
	Ratio of right d	ependency arcs 51	Con- and subjunction INCSC
TO COMPLEX THE STATE OF THE STA	Ratio of left de	pendency arcs 52	Past participles to verbs
	2 Modifier		
	3 Pre-modifi	er INCSC 54	Past verbs to verbs
	1 Post-modif	ier IncSc 55	Present verbs to verbs
	Subordina Subordina	te INCSC 56	Supine verbs to verbs
	Relative cla	A 1 (1 (1 (1 (1 (1 (1 (1 (1 (1	Relative structure INCSC
	7 Prepositional cor		Bilog type-token ratio
The same property of the same state of	Seman	Carlo Contract to the contract of the contract	하는 그 그 그 그 그 사람이 하면 하는 아이들이 가지 않는데 사람이 얼마를 하는데 하는데 하는데 되었다.
	3 Avg. nr. of sen	ses per token 60	이 그 그 아이들은 생겨보다 이번 이번 생활하다면 하다 생각하다면 하는데 살아내려면 하는데 없다.
· · · · · · · · · · · · · · · · · · ·	Noun sense	per noun 61	Pronouns to prepositions

Length-based features

- 1. Длина предложений
- 2. Средняя длина токена
- 3. Длинные слова (длиннее 13 символов)
- 4. Количество символов
- 5. LIX

Lexical features

KELLY Project (KEywords for Language Learning for Young and adults alike)

Вместо процентов – IncSc (incidence score) на 1000 слов: "The IncSc of a category was computed as 1000 divided by the number of tokens in the text or sentence multiplied by the count of the category in the sentence".

Lexical features

- 1. A1 lemma IncSc
- 2. A2 lemma IncSc
- 3. B1 lemma IncSc
- 4. B2 lemma IncSc
- 5. C1 lemma IncSc
- 6. C2 lemma IncSc
- 7. Difficult word IncSc (выше уровня текста)
- 8. Difficult noun and verb IncSc (выше уровня текста)

- 9. Out-of-Kelly IncSc (нет в списке)
 10. Missing lemma form IncSc (лемматизатор не распознал)
- 11. Avg. Kelly log frequency

Morphological features

- 1. Modal verbs to verbs
- 2. Particle IncSc
- 3. 3SG pronoun IncSc
- 4. Punctuation IncSc
- 5. Subjunction IncSc
- 6. S-verb IncSc
- 7. S-verbs to verbs
- 8. Adjective IncSc
- 9. Adjective variation (отношение категории к остальным категориям (сущ., прил., глаг., нареч.)

- 10. Adverb IncSc
- 11. Adverb variation
- 12. Noun IncSc
- 13. Noun variation
- 14. Verb IncSc
- 15. Verb variation
- 16. Nominal ratio
- 17. Nouns to verbs
- 18. Function word IncSc

Morphological features

- 19. Lexical words to nonlexical words (отношение к другим категориям)
- 20. Lexical words to all tokens
- 21. Neuter gender noun IncSc
- 22. Con- and subjunction IncSc
- 23. Past participles to verbs
- 24. Present participles to verbs
- 25. Past verbs to verbs
- 26. Present verbs to verbs
- 27. Supine verbs to verbs

- 28. Relative structure IncSc (consisted of relative adverbs, determiners, pronouns and possessives)
- 29. Bilog type-token ratio (Log all types / Log all token)
- 30. Square root type-token ratio (all types / root of all token),
- 31. Pronouns to nouns
- 32. Pronouns to prepositions

Syntactic features

- Average dependency length (MaltParser)
- Dependency arcs longer than 5
- 3. Longest dependency from root node
- Ratio of right dependency arcs
- Ratio of left dependency arcs
- 6. Modifier variation

- 7. Pre-modifier IncSc (e.g. adjectives and prepositional phrases)
- 8. Post-modifier IncSc (see 7)
- 9. Subordinate IncSc
- 10. Relative clause IncSc (E.g.: It is John (whom) Jack is waiting for.)
- 11. Prepositional complement IncSc

Semantic features

- 1. Avg. nr. of senses per token (использовали SALDO шведская альтернатива Princeton WordNet)
- 2. Noun senses per noun

Использовали WEKA (Waikato Environment for Knowledge Analysis), параметры по умолчанию (на Java)

- 1. a multinomial logistic regression model with ridge estimator (polytomous LR, multiclass LR, softmax regression, multinomial logit, the maximum entropy (MaxEnt) classifier, conditional maximum entropy model)
- 2. a multilayer perceptron,
- support vector machine learner, Sequential Minimal Optimization (SMO),
- 4. decision tree (J48).

Type	Number	Accuracy %	F	RMSE
Majority	-	33.2	0.17	0.52
LIX	1	34.9	0.22	0.38
<u>Lex</u>	<u>11</u>	<u>80.3</u>	<u>0.80</u>	<u>0.24</u>
<u>All</u>	<u>61</u>	<u>81.3</u>	<u>0.81</u>	<u>0.27</u>

multinomial logistic regression model with ridge estimator

Accuracy %

Type	Number	Perceptron	SM0	J48
Lex	11	<u>77.4</u>	42.1	55
AII	61	62.2	52.7	50.5

Confusion matrix

Predictions						
A 1	A2	B1	B2	C1		Label
37	12	0	0	0	A 1	
12	121	18	5	1	A2	
4	11	206	24	13	B1	
0	5	21	238	24	B2	
0	0	0	12	103	C 1	