

**Exact Methods in the Study of
Language and Text**

**Dedicated to Professor Gabriel Altmann
On the Occasion of His 75th Birthday**

**Edited by
Peter Grzybek & Reinhard Köhler**

**Mouton de Gruyter
Berlin – New York**

Contents

Viribus Quantitatis <i>Peter Grzybek and Reinhard Köhler</i>	v
A diachronic study of the style of Longfellow <i>Sergej N. Andreev</i>	1
Zum Gebrauch des deutschen Identitätspronomens ‘derselbe’ als funktionelles Äquivalent von Demonstrativ- und Personalpronomina aus historischer Sicht <i>John Ole Askedal</i>	13
Diversifikation bei Eigennamen <i>Karl-Heinz Best</i>	21
Bemerkungen zu den Formen des Namens <i>Schmidt</i> <i>Hermann Bluhme</i>	33
Statistical parameters of Ivan Franko’s novel <i>Perekhresni stežky (The Cross-Paths)</i> <i>Solomija Buk and Andrij Rovenchak</i>	39
Some remarks on the generalized Hermite and generalized Gegenbauer probability distributions and their applications <i>Mario Cortina-Borja</i>	49
New approaches to cluster analysis of typological indices <i>Michael Cysouw</i>	61
Menzerath’s law for the smallest grammars <i>Łukasz Dębowski</i>	77

Romanian online dialect atlas: Data capture and presentation <i>Sheila Embleton, Dorin Uritescu, and Eric Wheeler</i>	87
Die Ausdrucksmittel des Aspekts der tschechischen Verben <i>Jeehyeon Eom</i>	97
Quantifying the MULTEXT-East morphosyntactic resources <i>Tomaž Erjavec</i>	111
A corpus based quantitative study on the change of TTR, word length and sentence length of the English language <i>Fan Fengxiang</i>	123
On the universality of Zipf's law for word frequencies <i>Ramon Ferrer i Cancho</i>	131
Die Morrissche und die Bühlersche Triade – Probleme und Lösungs- vorschläge <i>Udo L. Figge</i>	141
Die kognitive Semantik der 'Wahrheit' <i>Michael Fleischer, Michał Grech, and Agnieszka Książek</i>	153
Kurzvorstellung der Korrelativen Dialektometrie <i>Hans Goebel</i>	165
A note on a systems theoretical model of usage <i>Johannes Gordesch and Peter Kunsmann</i>	181
Itemanalysen und Skalenkonstruktion in der Sprichwortforschung <i>Rüdiger Grotjahn und Peter Grzybek</i>	193
Do we have problems with Arens' law? A new look at the sentence- word relation <i>Peter Grzybek and Ernst Stadlober</i>	205
A language of thoughts is no longer an utopia <i>Wolfgang Hilberg</i>	219

Language subgrouping <i>Hans J. Holm</i>	225
Contextual word prominence <i>Luděk Hřebíček</i>	237
Das Menzerath-Gesetz in der <i>Vulgata</i> <i>Marc Hug</i>	245
Toward a theory of syntax and persuasive communication <i>Julian Jamison</i>	259
Grapheme und Laute des Russischen: Zwei Ebenen – ein Häufigkeitsmodell? Re-Analyse einer Untersuchung von A.M. Peškovskij <i>Emmerich Kelih</i>	269
Zur Zeitoptimierung der russischen Verbmorphologie <i>Sebastian Kempgen</i>	281
Ākāsha: between sphere and arrow – on the triple source for everything <i>Walter A. Koch</i>	287
Quantitative analysis of co-reference structures in texts <i>Reinhard Köhler and Sven Naumann</i>	317
Anthroponym – Pseudonym – Kryptonim: Zur Namensgebung in Erpresserschreiben <i>Helle Körner</i>	331
Quantitative linguistics within Czech contexts <i>Jan Králík</i>	343
Semantic components and metaphorization <i>Viktor Krupa</i>	353
Wortlängenhäufigkeit in J.W. v. Goethes Gedichten <i>Ina Kühner</i>	361

A general purpose ranking variable with applications to various ranking laws <i>Daniel Lavalette</i>	371
Wie schreibe ich einen Beitrag zu Gabriels Festschrift? <i>Werner Lehfeldt und [Lösung im Text]</i>	383
Bemerkungen zum Menzerath-Altmannschen Gesetz <i>Edda Leopold</i>	391
Die Stärkemessung des Zusammenhangs zwischen den Komponenten der Phraseologismen <i>Viktor Levickij and Iryna Zadorožna</i>	399
Pairs of corresponding discrete and continuous distributions: Mathematics behind, algorithms and generalizations <i>Ján Mačutek</i>	407
Linguistic numerology <i>Grigorij Ja. Martynenko</i>	415
Towards the measurement of nominal phrase grammaticality: contrasting definite-possessive phrases with definite phrases of 13 th to 19 th century Spanish <i>Alfonso Medina-Urrea</i>	427
A network perspective on intertextuality <i>Alexander Mehler</i>	439
Two semi-mathematical asides on Menzerath-Altman's law <i>Peter Meyer</i>	449
Stylometric experiments in modern Greek: Investigating authorship in homogeneous newswire texts <i>George K. Mikros</i>	461
On script complexity and the Oriya script <i>Panchanan Mohanty</i>	473

Statistical analogs in DNA sequences and Tamil language texts: rank frequency distribution of symbols and their application to evolutionary genetics and historical linguistics	485
<i>Sundaresan Naranan and Vriddhachalam K. Balasubrahmanyam</i>	
Zur Diversifikation des Bedeutungsfeldes slowakischer verbaler Präfixe	499
<i>Emília Nemcová</i>	
Ord's criterion with word length spectra for the discrimination of texts, music and computer programs	509
<i>Michael P. Oakes</i>	
Indexes of lexical richness can be estimated consistently with knowledge of elasticities: some theoretical and empirical results	521
<i>Epaminondas E. Panas</i>	
Huffman coding trees and the quantitative structure of lexical fields	533
<i>Adam Pawłowski</i>	
Linguistic disorders and pathologies: synergetic aspects	545
<i>Rajmund G. Piotrowski and Dmitrij L. Spivak</i>	
Text ranking by the weight of highly frequent words	555
<i>Ioan-Iovitz Popescu</i>	
Frequency analysis of grammemes vs. lexemes in Taiwanese	567
<i>Regina Pustet</i>	
Are word senses reflected in the distribution of words in text?	575
<i>Reinhard Rapp</i>	
Humanities' tears	587
<i>Jeff Robbins</i>	
Wortlänge im Polnischen in diachroner Sicht	597
<i>Otto A. Rottmann</i>	

The Menzerath-Altmann law in translated texts as compared to the original texts <i>Maria Roukk</i>	605
Different translations of one original text in a qualitative and quantitative perspective <i>Irma Sorvali</i>	611
The effects of diversification and unification on the inflectional paradigms of German nouns <i>Petra Steiner and Claudia Priin</i>	623
Nicht ganz ohne ... <i>Thomas Stolz, Cornelia Stroh and Aina Urdze</i>	633
Satz: stoisches axíōma oder peripatetischer lógos? <i>Wolf Thümmel</i>	647
Using Altmann-fitter for text analysis: An example from Czech <i>Ludmila Uhlířová</i>	659
Local grammars in word counting <i>Duško Vitas and Cvetana Krstev</i>	665
Fitting the development of periphrastic <i>do</i> in all sentence types <i>Relja Vulcanović and Harald Baayen</i>	679
Language change in a communication network <i>Eric S. Wheeler</i>	689
Die Suche nach Invarianten und Harmonien im Bereich symbolischer Formen <i>Wolfgang Wildgen</i>	699
Applying an evenness index in quantitative studies of language and culture: a case study of women's shoe styles in contemporary Russia <i>Andrew Wilson and Olga Mudraya</i>	709

The weighted mid-P confidence interval for the difference of independent binomial proportions <i>Viktor Witkovský and Gejza Wimmer</i>	723
Gabriel Altmann: Complete bibliography of scholarly works (1960–2005)	735
Tabula Gratulatoria <i>In Honor of Gabriel Altmann</i>	755

Do we have problems with Arens' law? A new look at the sentence-word relation

Peter Grzybek and Ernst Stadlober

Arens' Law owes its name to Gabriel Altmann who, in 1983, discussed the results of a book entitled *Verborgene Ordnung*, written by Hans Arens in 1965. In his book, Arens analyzed the specific relation between word length and sentence length; in detail, 117 samples of German literary prose texts were analyzed, written by 52 different authors. As a result, Arens arrived at the conclusion that an increase in sentence length goes along with an increase in word length. The raw data supporting this assumption can be reconstructed on the basis of the information given in Arens' book and are represented in Table 1. Calculating arithmetical means of word and sentence length (\bar{y} and \bar{x}), Arens presented his results in a graphical form, which implied a linear increase – cf. Figure 1a, p. 208. Two decades later, Altmann (1983) went a different way: in his discussion of Arens' findings, Altmann interpreted the observed relation in more general terms according to which the length of a particular (linguistic) component is a function of the length of the (linguistic) construct which it constitutes. This specific relation, which is well-known as Menzerath's Law today, was discussed by Altmann only a few years prior to his research on Arens' data. In his seminal "Prolegomena on Menzerath's Law", Altmann (1980) had suggested formula (1a) to be the most general form of what has hence been accepted to be the Menzerath-Altmann Law:

$$y = Ax^b e^{-cx}. \quad (1a)$$

In this context, Altmann had also presented two special cases of equation (1a), namely, equation (1b) for $c = 0$, and equation (1c) for $b = 0$.

$$y = Ax^b \quad (1b)$$

$$y = Ae^{-cx} \quad (1c)$$

Whereas equation (1a) is the most general form, equation (1b) has turned out to be the most commonly used "standard form" for linguistic purposes.

Table 1: Mean values for sentence length (\bar{x}) and word length (\bar{y}) for Arens' (1965) data, n denoting sample size in the number of words per sample

n	\bar{x}	\bar{y}	n	\bar{x}	\bar{y}	n	\bar{x}	\bar{y}
350	8.72	1.471	245	20.51	1.754	191	27.32	1.736
286	8.93	1.482	150	20.63	1.655	202	28.13	1.751
357	9.47	1.543	152	20.89	1.677	129	28.20	1.746
312	11.16	1.579	166	21.08	1.708	67	28.45	1.733
306	11.40	1.582	171	21.73	1.800	214	28.80	1.838
263	11.42	1.573	107	21.99	1.692	265	28.90	1.777
245	12.96	1.705	169	22.18	1.689	103	29.39	1.789
131	13.36	1.596	205	22.44	1.717	105	29.50	1.737
249	13.50	1.591	133	22.62	1.829	131	29.81	1.813
478	13.65	1.662	210	22.66	1.716	116	30.65	1.774
388	13.66	1.603	132	22.74	1.691	137	30.70	1.775
223	13.84	1.602	479	23.14	1.658	140	30.80	1.771
290	13.92	1.613	160	23.48	1.692	204	30.93	1.806
575	14.07	1.683	399	23.52	1.723	120	31.03	1.777
213	14.13	1.649	247	24.15	1.739	139	31.34	1.820
276	14.53	1.670	129	24.22	1.737	145	31.14	1.780
302	14.70	1.617	124	24.27	1.759	97	32.67	1.752
397	15.13	1.593	200	24.31	1.709	93	32.84	1.794
205	15.40	1.651	124	24.33	1.727	88	34.06	1.799
256	15.60	1.668	123	24.48	1.729	95	34.11	1.801
389	15.85	1.733	218	24.50	1.714	122	34.84	1.763
451	16.23	1.628	200	24.70	1.711	206	35.32	1.762
200	16.37	1.628	272	24.90	1.580	87	35.41	1.727
363	16.53	1.631	166	25.00	1.698	141	35.95	1.945
257	16.57	1.777	154	25.07	1.717	100	36.02	1.779
254	16.73	1.676	211	25.10	1.673	225	36.52	1.722
181	16.91	1.764	166	25.13	1.814	82	37.52	1.761
200	17.22	1.639	119	25.27	1.725	148	37.61	1.777
202	17.23	1.635	118	25.42	1.721	301	37.94	1.842
210	17.65	1.664	110	25.53	1.724	122	38.17	1.851
191	18.37	1.660	125	26.00	1.727	78	39.23	1.863
407	19.68	1.683	135	26.02	1.755	81	39.67	1.847
223	19.69	1.711	334	26.07	1.600	82	40.29	1.830
158	19.70	1.661	200	26.35	1.784	84	41.20	1.871
243	19.98	1.682	160	26.40	1.827	124	42.65	1.805
230	20.00	1.678	212	27.00	1.752	100	42.74	1.895
200	20.02	1.678	255	27.19	1.739	148	45.41	1.819
200	20.05	1.670	176	27.19	1.713	70	60.76	1.817
229	20.14	1.782	150	27.30	1.699	73	92.40	1.935

With regard to the relation between sentence length and word length, Altmann (1983: 31) pointed out that Menzerath's Law as described above is likely to hold true only when one is concerned with the direct constituents of a given construct. Therefore, in its direct form, Menzerath's Law might fail to grasp the relation between sentence length and word length, as soon as we are not concerned with the word as the direct constituent of the sentence.

In fact, an intermediate level is likely to come into play – such as for example phrases or clauses as the direct constituents of the sentence. In this case, words might well be the direct constituents of clauses or phrases, but they would only be indirect constituents of a sentence. Consequently, an increase in sentence length should result in an increase in word length, too — as in fact observed by Arens. Corresponding observations must therefore not be misinterpreted in terms of a counterproof to Menzerath's Law; rather, they should be understood as an indirect proof of it in the form of Arens' Law. Yet, according to Arens's Law, as described by Altmann, the increase in word length with increasing sentence length should not be linear; rather it should follow Menzerath's Law. Strictly speaking, with y symbolizing word length, z symbolizing phrase (or clause) length, and x symbolizing sentence length, we were thus concerned with two relations simultaneously: $y = Az^b e^{cz}$ and $z = A'x^{b'} e^{c'x}$. Inserting the latter equation into the first, one obtains y as a function of x :

$$y = A''x^{b''} e^{c''x + A'''x^{b'} e^{c'x}} . \quad (2)$$

Given that the “standard case” of Menzerath's Law (1b) has often been sufficient to describe the relation between sentence length and clause length (i.e., $z = Ax^b$), as well as the one between clause length and word length (i.e., $y = A'z^{b'}$), Altmann (1983: 32) argued in favor of using this special case, consequently obtaining $y = A''x^{b''}$, corresponding to equation (1b). The only difference to be expected for the relation between directly and indirectly related units of different levels is that, in case of directly neighboring units, parameters b and b' should be negative (due to the prognosed decline); in case of indirectly related units, with intermediate levels, $b'' = b \cdot b'$ will become positive. In addition to the linear regression, Figure 1b represents the results for fitting equations (1a) and (1b) to Arens' data.

Testing the goodness of fit of the non-linear Menzerathian model (1b) with $\hat{y} = 1.2183x^{0.1089}$, Altmann calculated an F -test which, with $\hat{F}_{1,115} = 241.40$, he interpreted to be a highly significant result, corroborating his assumptions on the Menzerathian relation between sentence length and word length. This

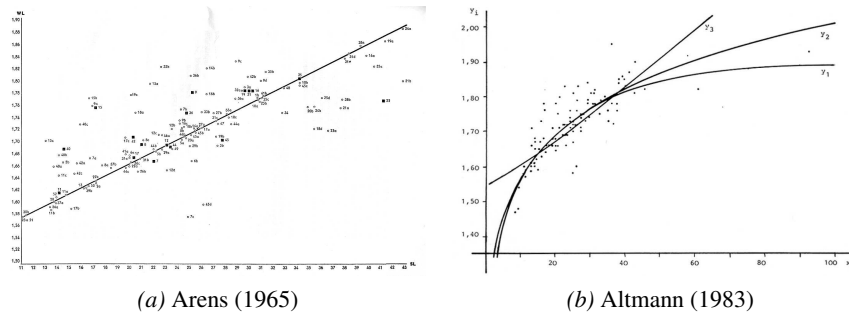


Figure 1: Sentence length and word length: linear and non-linear regression

regularity describing the dependence of units from two indirectly related linguistic levels has henceforth become well-known by the name of Arens' Law (or Arens-Altman Law). Yet, taking a second look at Altmann's (1983) modeling of Arens' data, doubt may arise with regard to two points, and they even give rise to the fundamental question whether we have problems with Altmann-Arens' Law:

1. First, a decade after Altmann's (1983) study, Grotjahn (1992) discussed some methodological weaknesses of the F -test for testing linguistic data; as a result, Grotjahn argued in favor of calculating the determination coefficient R^2 , instead of F -tests, favorably in form of equation (3).

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}. \quad (3)$$

Now, re-analyzing Arens' data according to equation (3), results in a rather poor value of $R^2 = 0.70$ (a value of $R^2 \geq 0.85$ usually being assumed to indicate a satisfying fit). Thus, notwithstanding the fact that the result for the non-linear regression model is definitely better than the one for the linear model (with $R^2 = 0.58$), it is far from being convincing, consequently shedding doubt on the adequacy of the Menzerathian interpretation.

2. Second, the scope of Menzerath's Law initially has been to describe the relation between the constituting components of a given construct; consequently, Menzerath's Law must be understood as having been designed in terms of an **intra-textual** law, relevant for the internal struc-

ture of a given text sample.¹ Arens' data, however, are of a different kind, implying **inter-textual** relations, based on the calculation of the mean lengths of words (\bar{x}_i) and sentences (\bar{y}_i) for each of the 117 text samples, resulting in two vectors of arithmetic means (\bar{x}_i and \bar{y}_i).

Altmann (1983: 32), who based his analyses on these vectors, was of course well aware of the difference between intra- and inter-textual dependences (though not explicitly using these words), and he emphasized that Arens' data cannot be taken as a direct proof of the Menzerathian relation on an intra-textual level. Still, he interpreted Arens' (inter-textual) data to be even more reliable, likely to rule out possibly intervening (intra-textual) individual variances. Yet, principally speaking, it must be noted that we are concerned with two different applications, or interpretations, of what has been discussed as Arens' Law:

- in an intra-textual perspective, Arens' Law may be interpreted to be a logical derivation of Menzerath's Law, due to the intervention of intermediate levels (cf. Altmann & Schwibbe 1989: 12f., Cramer 2005);
- in an inter-textual perspective, Arens' Law is not necessarily a logical consequence of Menzerath's Law; rather, it has the very same status of a strong hypothesis as has Menzerath's Law itself.²

In summary, we are thus faced with two possibly interrelated problems which ask for clarification:

1. interpreting the relation between sentence length and word length along the Altmann-Arens line, one must separate the intra-textual and inter-textual implications more clearly than this has been done hitherto;
2. the poor empirical evidence in support of the Altmann-Arens Law outlined above gives rise to the question of possible reasons for this circumstance.

1. We need not discuss the notion of 'text' here; for the sake of simplification we tolerate that a 'text' may be represented by homogeneous material, as well as by a mixed corpus, or by dictionary material, etc.

2. Given Arens' Law is relevant on the intra-textual level, this is no indication of a general increase in word length with an increase in sentence length, on the inter-textual level: With regard to the intra-textual level, Arens' Law means that the mean word length is an increasing function of sentence length. In comparing texts on an inter-textual level, we take only mean word length and mean sentence length of each text and study the relationship between these means across different texts; we cannot suspect the same rule applies as on the intra-textual level.

The present text concentrates on the inter-textual perspective, and it focusses on possible explanations for the obviously poor results in the context of Arens' data. It seems reasonable to start from this inter-textual end, tentatively maintaining Altmann's (1983: 32) assumption as to less variance across samples than for individual texts, consequently predicting even worse results for individual texts (i.e., for the intra-textual situation). A clarification of the inter-textual level might therefore provide important insight into the mechanism of Arens' Law, in general, and should thus yield valuable results for future intra-textual studies (cf. Grzybek et al. 2006). As to the observed poverty of the results, it seems important to take into account the circumstance that Arens' Law, as well as Menzerath's Law, has been designed as what one might term a 'law of averages'. This is to say that the application of these laws to linguistic data has been guided by the interest to express overall tendencies within larger linguistic samples: to this end, arithmetical means have been calculated for particular data points, and the means of particular independent variables (\bar{x}) have been related to the means of the relevant dependent variables (\bar{y}). In case of the relation between sentence length and word length, we are concerned with two arithmetical means: \bar{x}_i as the independent variable denoting average sentence length, and \bar{y}_i as the dependent variable denoting the corresponding word length. As was mentioned above, in case of inter-textual studies, we thus obtain two vectors of arithmetic means, $\bar{\mathbf{x}}_i$ and $\bar{\mathbf{y}}_i$; in this case, for $i = 1 \dots N$ texts, each individual average value \bar{x}_i and \bar{y}_i is based on a particular number of observations within the text. Yet, due to the large variance of sentence length and the resulting great amount of classes³, we tend to have only one single \bar{y}_i value for each data point⁴ of

3. This is the reason why "simple" sentence length studies, focusing on mere frequency distributions of sentence length, tend to form particular intervals (usually of five classes), rather than take into consideration each individual sentence length class.

4. The situation may be less complicated when applying Arens' Law to other linguistic levels, on which the number of linguistic classes is limited, in practice. This is particularly evident in case of Menzerathian studies of word length (in terms of the number of syllables, or of morphemes, per word); but also for level-transgressing studies (implying Arens' Law), when word length is measured by the number of letters or phonemes per word, the number of classes still is small. As opposed to this, in case of sentence length, the variation is much larger; this is less relevant for Menzerathian studies (measuring sentence length by the number of clauses per sentence) than it is for Arens studies (based on the number of words per sentence): as a consequence, it is a mere fact of coincidence that two texts (albeit only two) have an identical average sentence length. In fact, as an inspection of the 117 texts represented in Table 1 shows, this occurs only once in Arens' data (namely, for $\bar{x} = 27.19$).

the independent variable \bar{x}_i . The interpretation of data in terms of Arens' Law may therefore be deluded by the fact that, although the averages are based on a rather large number of observations, for each independent data point \bar{x}_i being introduced into the regression model, there is only a single dependent value (\bar{y}_i). It seems to be reasonable therefore to test in how far some kind of **data pooling**, providing some kind of "second-order" averages, will lead to more satisfying results. However, pooling itself is not unproblematic, the more since there are different pooling procedures:

1. either one defines a particular (minimal) *number of observations* for calculating the mean value.
2. or one considers all data points within a given *interval* and calculates the corresponding arithmetical mean;

Both procedures imply a certain degree of subjective arbitrariness, since neither the concrete number of observations nor the interval size can be theoretically defined a priori. And even having made a decision for a particular interval size, the next problem which arises concerns the lower limit of the first interval: given a desired interval of five, for example: should the first interval start with 1 (a theoretical minimum), with 2 (one possibility to linguistically justify a sentence length minimum), or with 8 (the observed minimum in the given sample)? Obviously, there can only be an authoritative decision – favoring an empirically based optimum would cause variations from one sample to another (and, consequently complicate between-text comparisons). Additionally, results are likely to be influenced by the decision to calculate either 'simple' arithmetical means or weighted means (thus taking into consideration the number of sentences on which the observation is based).

In the context of these factors asking for a decision, a number of logical consequences must not be ignored which are of utmost importance. Thus, if we decide to have rather large classes or intervals (in order to have more observations within a given class), one must be aware of the fact that this will result in fewer data points making the interpretation more difficult (unless one has an abundant mass of data). Smaller groups or intervals, however, will lead to the fact that many data points may be represented on the basis of a relatively narrow segment of the whole data spectrum. Thus, not only is there no pooling procedure which may be favored for theoretical reasons; additionally, none of these procedures is unproblematic in practice. In fact, any decision made is likely to be a secondary factor influencing the result, which may be highly dependent on the specific data structure under study.

The aim of the present study is not so much to offer solutions to all open questions, as to point out general problems in dealing with Arens' Law, which are, among others, related to the problem of pooling. Let us therefore, by way of an example, re-analyse Arens' original data (cf. Table 1, p. 206). Table 2 represents the pooled data, each data class based on five observations, the original data sorted in ascending order of sentence length (\bar{x}).⁵

Table 2: Mean values for sentence length (\bar{x}) and word length (\bar{y}) for Arens' (1965) data, in classes of five observations

i	f	\bar{x}	\bar{y}	\hat{y}	i	f	\bar{x}	\bar{y}	\hat{y}
1	1–5	9.936	1.531	1.568	13	60–65	24.954	1.676	1.731
2	6–10	12.978	1.625	1.614	14	66–70	25.470	1.742	1.734
3	11–15	13.924	1.630	1.626	15	71–75	26.368	1.744	1.741
4	16–20	15.072	1.640	1.640	16	76–80	27.426	1.728	1.748
5	21–25	16.310	1.679	1.654	17	81–85	28.748	1.777	1.757
6	26–30	17.148	1.676	1.663	18	86–90	30.292	1.774	1.767
7	31–35	19.484	1.679	1.685	19	91–95	31.422	1.787	1.774
8	36–40	20.144	1.712	1.691	20	96–100	34.234	1.784	1.790
9	41–45	21.264	1.706	1.701	21	101–105	36.284	1.787	1.801
10	46–50	22.528	1.728	1.712	22	106–110	38.524	1.836	1.813
11	51–55	23.702	1.710	1.721	23	111–117	52.207	1.853	1.873
12	56–60	24.378	1.728	1.726					

Figure 2a illustrates the convincing result, characterized by a determination coefficient of $R^2 = 0.93$ for parameter values $a = 1.2268$ and $b = 0.1070$.

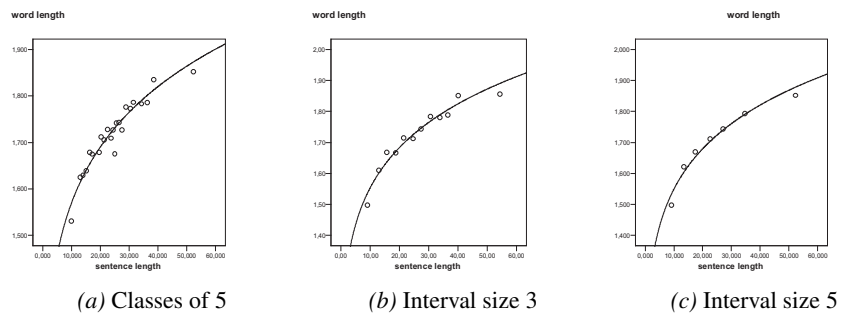


Figure 2: Arens' (1965) data, with different kinds of pooling

5. Given a sample size of 117 texts, the last class includes seven observations.

By way of a comparison, Table 3 gives the results of fitting equation (1b) to the data, pooled according to intervals, with two interval sizes: five and three; the number of observations the calculated mean is based on is indicated by n , the theoretical values are indicated by \hat{y} . As can be seen from Figures 2b

Table 3: Mean values for sentence length (\bar{x}) and word length (\bar{y}) for Arens' (1965) data, in intervals of length 3 vs. 5

i	n	Interv.	\bar{x}	\bar{y}	\hat{y}	i	n	Interv.	\bar{x}	\bar{y}	\hat{y}
1	3	[8,11)	9.040	1.499	1.535	1	3	[5,10)	9.040	1.499	1.532
2	10	[11,14)	12.887	1.611	1.600	2	14	[10,15)	13.307	1.623	1.602
3	14	[14,17)	15.625	1.669	1.636	3	28	[15,20)	17.269	1.671	1.651
4	8	[17,20)	18.690	1.667	1.670	4	27	[20,25)	22.581	1.712	1.704
5	15	[20,23)	21.312	1.716	1.696	5	25	[25,30)	26.982	1.744	1.739
6	19	[23,26)	24.554	1.713	1.724	6	13	[30,35)	34.542	1.795	1.790
7	15	[26,29)	27.288	1.745	1.745	7	7	35+	52.207	1.853	1.878
8	10	[29,32)	30.529	1.784	1.768						
9	5	[32,35)	33.704	1.782	1.789						
10	8	[35,38)	36.536	1.789	1.805						
11	4	[38,41)	40.098	1.853	1.825						
12	6	41+	54.193	1.857	1.890						

and 2c, the results for fitting equation (1b) to Arens' data are very convincing, irrespective of interval size:

1. For intervals of three, the determination coefficient is $R^2 = 0.95$ with parameter values $a = 1.1887$ and $b = 0.1161$ – cf. 2b.
2. The result is equally fine, when the means are based on intervals of five: in this case, the determination coefficient is $R^2 = 0.97$, with parameter values $a = 1.1856$ and $b = 0.1163$ – cf. Figure 2c.⁶

Data pooling thus in fact turns out to be a crucial matter in dealing with Arens' data and, consequently, with Arens' Law. If the first conclusion therefore is that proving Arens' Law demands some kind of data pooling in order for the overall tendency to become transparent, then the second conclusion implies the availability of sufficient data material when studying Arens' Law (at least on an inter-textual level).

6. A regression analysis which is not based on the a priori defined intervals given in Table 3, but – given a minimal sentence length of 8.72 –, starts with a lower interval border of 8 – thus including intervals of [8,13), [13,18), [18,23), ... –, leads to an almost identical result of $R^2 = 0.98$.

Yet, a large amount of data is a necessary, but not a sufficient condition. Rather, in dealing with Arens' Law, due attention must be paid to the factor of data homogeneity. This shall be demonstrated here by enlarging our data base of Arens' texts with relevant data presented by Wilhelm Fucks (1955, 1956) a decade before Arens' work. In his pioneering studies on the mathematics of literary style, Fucks studied the relation between sentence length and word length, though not concentrating on a mathematical model of this relation. Still, he provided relevant data of 54 German text samples;⁷ half of them were literary prose, the other half scholarly prose.

Combining Fucks' and Arens' data into one common corpus of 171 text samples, one might expect the result to improve as compared to Arens' data alone; yet, re-analyzing the relation between sentence length and word length of the joint corpus according to equation (1b), results in a very poor value of $R^2 = 0.22$, which is not only far from being satisfying, but, more importantly, significantly worse as compared to the result obtained above for Arens' data alone (with $R^2 = 0.70$).

Searching for a reason of this deterioration, it seems reasonable to follow Fucks' initial ideas assuming that the two groups of texts belong to two different writing styles, characterized by differences in sentence length and word length (cf. Table 4).

Table 4: Comparison of literary and scholarly prose (Fucks 1955 and Arens 1965)

		N	Word Length		Sentence Length	
			\bar{y}	s	\bar{x}	s
Arens	Literature	117	1.72	0.09	25.37	10.92
Fucks	Literature	27	1.68	0.09	19.28	5.61
Fucks	Prose	27	1.98	0.13	24.39	6.56

This can clearly be seen from Figure 3a (taken from Fucks 1955: 239), which shows that the two text groups are separated mainly along the vertical axis, the differences thus being related to differences in word length rather than sentence length. A re-analysis of Fucks' data by way of a discriminant analysis confirms this impression: only 61.10% of the texts are correctly classified with sentence length as the relevant discriminant variable, as compared

7. There is an important difference between Fucks' and Arens' data: whereas Arens' analyzed coherent text segments of at least 3 000 words, Fucks combined five randomly chosen segments of 500 words each (cf. Arens 1965: 16).

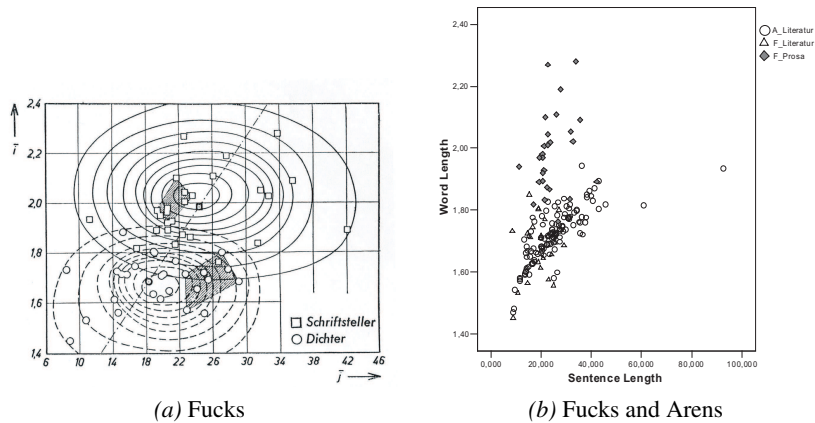


Figure 3: Sentence length and word length (Fucks 1955 and Arens 1965)

to 92.60% correct classifications on the basis of word length. This tendency is also reflected in the graphical representation of the combined corpus; as can be seen from Figure 3b, Fucks' literary prose texts neatly fit with the group of Arens's data, whereas the scholarly prose texts clearly fall into a different area.

As can easily be seen (and, in fact, statistically proven by way of post-hoc tests), the two literature samples fall into one category as to word length, but they differ significantly in sentence length (both as to \bar{x} and s); on the other hand, Arens' literary texts and Fucks' scholarly prose texts fall into one category as to sentence length (though with enormously differing s), but clearly differ in word length. Details as to possible reasons for this rather unexpected result need not be discussed here. It is well possible that the observed differences are partly related to the randomness of Fucks's data samples, or to diverging definitions of 'word' and/or 'sentence'. Yet, such (additional) factors are not likely to explain the whole complexity of the matter: and even if they should explain differences between the two samples of literary texts, the specifics of the scholarly prose texts make it most obvious that we are concerned with a specific group of texts. It seems unlikely, therefore, that all texts follow one common tendency. This conclusion is of utmost importance for the relevance of Arens' Law, with regard to which we have to conclude that, on an inter-textual level, it is likely to be operative only within homogeneous text groups, if at all. In fact, it may well turn out that, as soon as we concen-

trate on homogeneous groups of texts only, the latter do not display enough variance of either word or sentence length, due to genre specific structures. If this were true, Arens' Law were not likely to become transparent within a given text group and, on the inter-textual level, would at best turn out by way of a text type related law. In fact, of our three samples, only Arens' literary texts vary sufficiently with regard to both sentence and word length; here, pooling turns out to be a necessary and efficient procedure. As compared to this, analyzing the two Fucks samples (literary vs. scholarly prose) separately, not only results in extremely poor values of ($R^2 = 0.07$) and ($R^2 = 0.10$), respectively; additionally, in this case, pooling makes no sense due to the small sample sizes of $N = 27$. The question must remain open for further research (cf. Grzybek et al. 2006), therefore, what will happen to the assumptions suggested by Arens' Law as soon as one analyzes sufficient homogeneous data.

In summary, possible problems with Arens's Law may be related to different factors:

1. Attention must be paid to the distinction of intra-textual and inter-textual perspectives when dealing with Arens' Law.
2. It seems reasonable that Arens' Law is valid only within the framework of particular text sorts, or discourse types;
3. Arens's Law seems to express specific tendencies which can be submitted to observation only in case of large data material, or by way of specific pooling procedures; pooling, in turn, may lead to partly diverging results, depending on the concrete procedure chosen.

References

- Altmann, Gabriel
 1980 "Prolegomena to Menzerath's law". In: *Glottometrika 2*. Bochum: Brockmeyer, 1–10.
- 1983 "H. Arens' «Verborgene Ordnung» und das Menzerathsche Gesetz". In: Faust, Manfred; Harweg, Roland; Lehfeldt, Werner; Wienold, Götz (Hg.), *Allgemeine Sprachwissenschaft, Sprachtypologie und Textlinguistik*. Tübingen: Narr, 31–39.
- Altmann, Gabriel; Schwibbe, Michael H.
 1989 *Das Menzerathsche Gesetz in informationsverarbeitenden Systemen*. Hildesheim: Olms.

- Arens, Hans
1965 *Verborgene Ordnung. Die Beziehungen zwischen Satzlänge und Wortlänge in deutscher Erzählprosa vom Barock bis heute.* Düsseldorf: Pädagogischer Verlag Schwann.
- Cramer, Irene M.
2005 "Das Menzerathsche Gesetz". In: Köhler, Reinhard; Altmann, Gabriel; Piotrowski, Raimund G. (Eds.), *Quantitative Linguistik. Ein internationales Handbuch.* Berlin / New York: de Gruyter, 659–688.
- Fucks, Wilhelm
1955 "Unterschied des Prosastils von Dichtern und Schriftstellern. Ein Beispiel mathematischer Stilanalyse." In: *Sprachforum*, 1; 234–241.
- Grotjahn, Rüdiger
1993 "Evaluating the adequacy of regression models: some potential pitfalls". In: *Glottometrika 13.* Bochum: Brockmeyer, 121–172.
- Grzybek, Peter; Kelih, Emmerich; Stadlober, Ernst
2006 "The relationship of word length and sentence length: the inter-textual perspective" In: *Advances in Data Analysis.* Heidelberg /New York: Springer. [In print]