

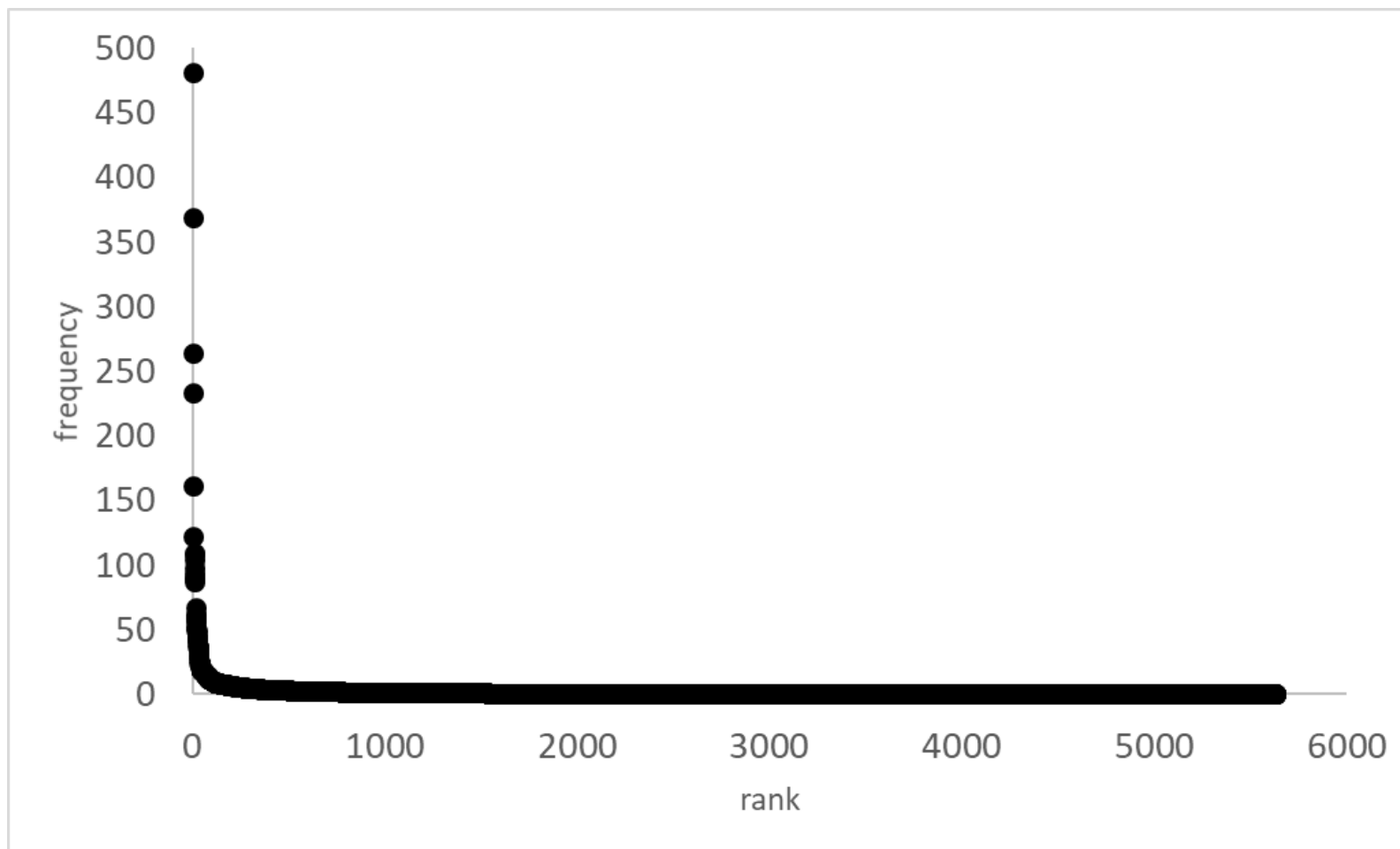
Stylistika VII

ZS 2024

Obsah

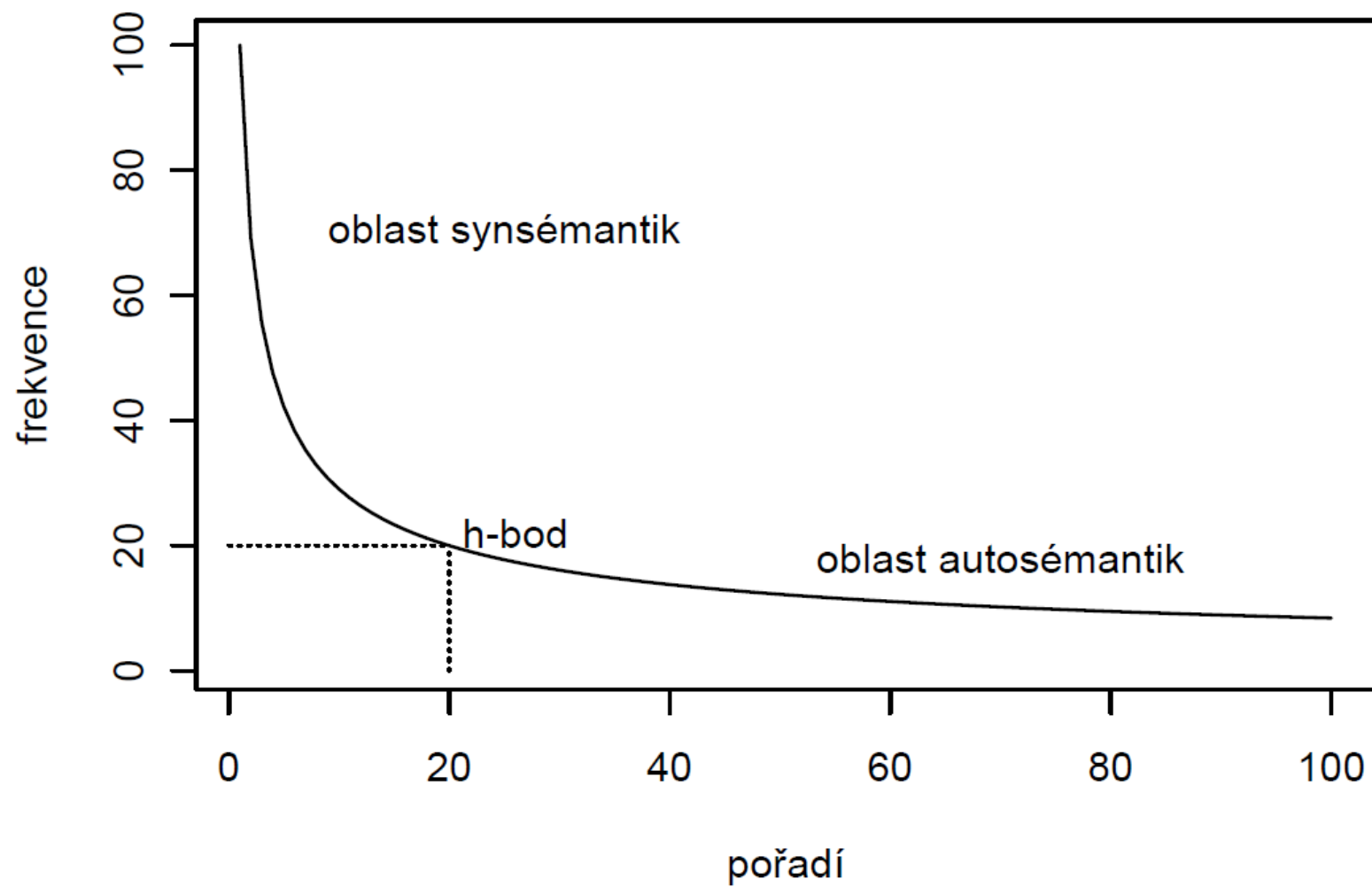
- klasifikace na základě nejfrekventovanějších slov
- literatura
 - Plechač. P. (2020). Jak určit autora textu. Vesmír, 99, s. 210-211)
 - <https://vesmir.cz/cz/casopis/archiv-casopisu/2020/cislo-4/jak-urcit-autora-textu.html>

Frekvenční struktura textu



J. Škvorecký: Eva byla nahá

Frekvenční struktura textu



Nejfrekventovanější slova

Škvorecký: Eva byla nahá		
pořadí	slovo	f
1	a	481
2	se	369
3	na	264
4	v	234
5	jsem	161
6	s	122
7	z	110
8	američan	108
9	to	104
10	řekl	98
11	ale	94
12	do	91
13	že	89
14	řekla	87
15	dívka	67

Nejfrekventovanější slova

Škvorecký: Eva byla nahá		
pořadí	slovo	f
1	a	481
2	se	369
3	na	264
4	v	234
5	jsem	161
6	s	122
7	z	110
8	američan	108
9	to	104
10	řekl	98
11	ale	94
12	do	91
13	že	89
14	řekla	87
15	dívka	67

Nejfrekventovanější slova

Škvorecký: Eva byla nahá		
pořadí	slovo	f
1	a	481
2	se	369
3	na	264
4	v	234
5	jsem	161
6	s	122
7	z	110
8	američan	108
9	to	104
10	řekl	98
11	ale	94
12	do	91
13	že	89
14	řekla	87
15	dívka	67

Hrabal: Perlička na dně		
pořadí	slovo	f
1	a	2239
2	se	1203
3	to	1037
4	na	879
5	ale	514
6	tak	504
7	do	467
8	si	459
9	jsem	456
10	v	446
11	že	440
12	je	432
13	já	363
14	když	296
15	jak	283

Hašek: Osudy...I.		
pořadí	slovo	f
1	a	7045
2	se	6061
3	na	3927
4	že	3469
5	to	3075
6	v	2585
7	je	1801
8	do	1749
9	s	1667
10	si	1534
11	když	1387
12	z	1375
13	tak	1308
14	jsem	1286
15	švejk	1188

Nejfrekventovanější slova

Škvorecký: Eva byla nahá		
pořadí	slovo	f_rel
1	a	0.037
2	se	0.028
3	na	0.020
4	v	0.018
5	jsem	0.012
6	s	0.009
7	z	0.008
8	američan	0.008
9	to	0.008
10	řekl	0.007
11	ale	0.007
12	do	0.007
13	že	0.007
14	řekla	0.007
15	dívka	0.005

Hrabal: Perlička na dně		
pořadí	slovo	f_rel
1	a	0.054
2	se	0.029
3	to	0.025
4	na	0.021
5	ale	0.012
6	tak	0.012
7	do	0.011
8	si	0.011
9	jsem	0.011
10	v	0.011
11	že	0.011
12	je	0.010
13	já	0.009
14	když	0.007
15	jak	0.007

Hašek: Osudy...I.		
pořadí	slovo	f_rel
1	a	0.035
2	se	0.030
3	na	0.020
4	že	0.017
5	to	0.015
6	v	0.013
7	je	0.009
8	do	0.009
9	s	0.008
10	si	0.008
11	když	0.007
12	z	0.007
13	tak	0.007
14	jsem	0.006
15	švejk	0.006

Nejfrekventovanější slova

Škvorecký: Eva byla nahá		
pořadí	slovo	f_rel
1	a	0.037
2	se	0.028
3	na	0.020
4	v	0.018
5	jsem	0.012
6	s	0.009
7	z	0.008
8	američan	0.008
9	to	0.008
10	řekl	0.007
11	ale	0.007
12	do	0.007
13	že	0.007
14	řekla	0.007
15	dívka	0.005

Hrabal: Perlička na dně		
pořadí	slovo	f_rel
1	a	0.054
2	se	0.029
3	to	0.025
4	na	0.021
5	ale	0.012
6	tak	0.012
7	do	0.011
8	si	0.011
9	jsem	0.011
10	v	0.011
11	že	0.011
12	je	0.010
13	já	0.009
14	když	0.007
15	jak	0.007

Hašek: Osudy...I.		
pořadí	slovo	f_rel
1	a	0.035
2	se	0.030
3	na	0.020
4	že	0.017
5	to	0.015
6	v	0.013
7	je	0.009
8	do	0.009
9	s	0.008
10	si	0.008
11	když	0.007
12	z	0.007
13	tak	0.007
14	jsem	0.006
15	švejk	0.006

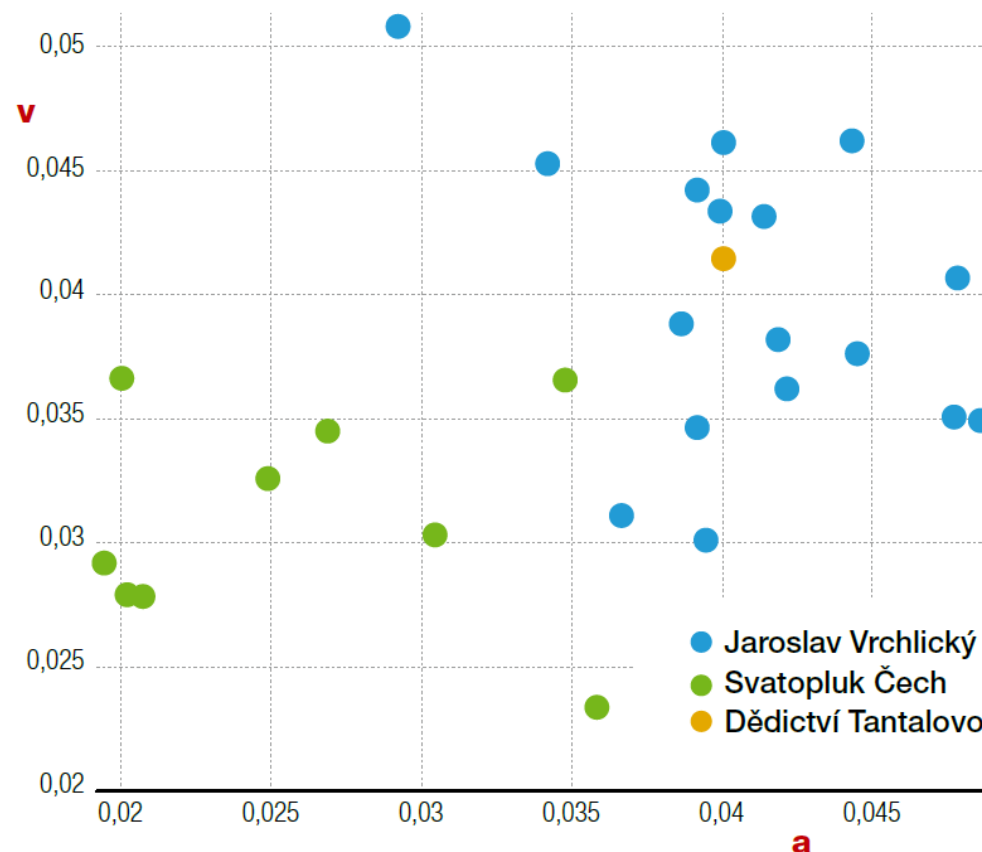
Nejfrekventovanější slova

Škvorecký: Eva byla nahá		
pořadí	slovo	f_rel
1	a	0.037
2	se	0.028
3	na	0.020
4	v	0.018
5	jsem	0.012
6	s	0.009
7	z	0.008
8	američan	0.008
9	to	0.008
10	řekl	0.007
11	ale	0.007
12	do	0.007
13	že	0.007
14	řekla	0.007
15	dívka	0.005

Hrabal: Perlička na dně		
pořadí	slovo	f_rel
1	a	0.054
2	se	0.029
3	to	0.025
4	na	0.021
5	ale	0.012
6	tak	0.012
7	do	0.011
8	si	0.011
9	jsem	0.011
10	v	0.011
11	že	0.011
12	je	0.010
13	já	0.009
14	když	0.007
15	jak	0.007

Hašek: Osudy...I.		
pořadí	slovo	f_rel
1	a	0.035
2	se	0.030
3	na	0.020
4	že	0.017
5	to	0.015
6	v	0.013
7	je	0.009
8	do	0.009
9	s	0.008
10	si	0.008
11	když	0.007
12	z	0.007
13	tak	0.007
14	jsem	0.006
15	švejk	0.006

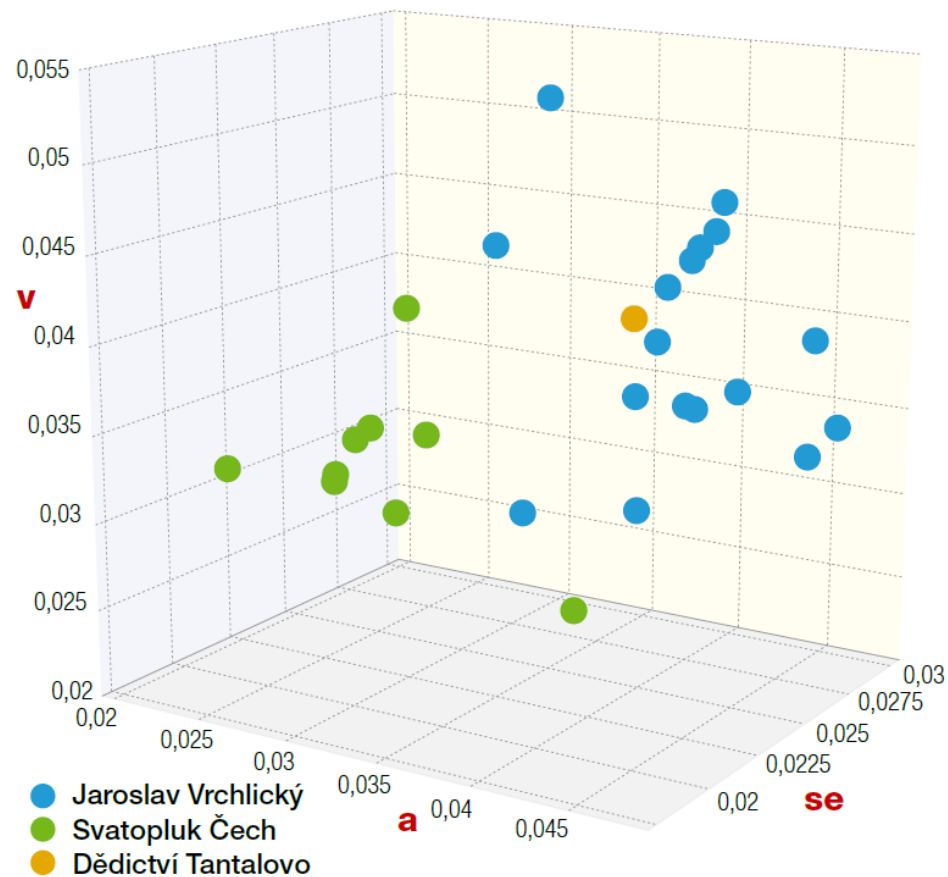
Nejfrekventovanější slova



**1. RELATIVNÍ
ČETNOSTI** dvou
nejfrekventova-
nějších slov
(*a*, *v*) ve sbírkách
Jaroslava
Vrchlického,
Svatopluka
Čecha a sbírce
*Dědictví
Tantalovo*.

Plechač (2020)

Nejfrekventovanější slova



2. RELATIVNÍ ČETNOSTI tří nejfrekventovanějších slov (a, v, se) ve sbírkách Jaroslava Vrchlického, Svatopluka Čecha a sbírce *Dědictví Tantalovo*.

Plechač (2020)

Nejfrekventovanější slova

- nezávislost na tématu
- aproximace gramatiky (viz níže)

Měření – Burrowsova delta

- relativní frekvence daného počtu nejfrekventovanějších slov
- výpočet z-skóre
- pro každou dvojici textů se pro každé slovo (z daného počtu nejfrekventovanějších slov) vypočítá rozdíl hodnot z-skóre
- tyto vzdálenosti se sečtou
- aplikuje se shluková analýza

Modelový příklad

- text A a text B
- vypočet na základě pouze dvou nejfrekventovanějších slov
 - jedná se jen o modelový příklad, obvykle se pracuje se stovkami slov

	a	na	délka textu
text A	56	32	800
text B	27	90	900

Modelový příklad

- relativní frekvence (v %)

$$f_{rel(text A)}(a) = \frac{f(a)}{N_{(text A)}} 100 = \frac{56}{800} 100 = 0.07 \cdot 100 = 7$$

	a	na	délka textu
text A	56	32	800
text B	27	90	900

Modelový příklad

- relativní frekvence (v %)

$$f_{rel(text A)}(a) = \frac{f(a)}{N_{(text A)}} 100 = \frac{56}{800} 100 = 0.07 \cdot 100 = 7$$

$$f_{rel(text B)}(a) = \frac{f(a)}{N_{(text B)}} 100 = \frac{27}{900} 100 = 0.03 \cdot 100 = 3$$

$$f_{rel(text B)}(na) = \frac{f(na)}{N_{(text B)}} 100 = \frac{32}{800} 100 = 0.04 \cdot 100 = 4$$

$$f_{rel(text A)}(na) = \frac{f(na)}{N_{(text B)}} 100 = \frac{90}{900} 100 = 0.1 \cdot 100 = 10$$

Modelový příklad

- z-skóre
 - **rozdíl** mezi **relativní frekvencí** daného slova v textu a **průměrnou relativní frekvencí** daného slova ve vzorku sestaveného z analyzovaných textů vydělený **směrodatnou odchylkou**

Modelový příklad

- z-skóre
 - **rozdíl** mezi **relativní frekvencí** daného slova v textu a **průměrnou relativní frekvencí** daného slova ve vzorku sestaveného z analyzovaných textů vydělený **směrodatnou odchylkou**
 - v tomto modelovém případě si představme, že máme více textů, ze kterých jsme spočítali průměrnou relativní frekvenci a směrodatnou odchylku

Variabilita dat – rozptyl & směrodatná odchylka

- rozptyl
 - střední hodnota kvadrátů odchylek od střední hodnoty

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_N - \bar{x})^2}{N - 1}$$

Variabilita dat – rozptyl & směrodatná odchylka

{2,2,3,3,4,5}

průměr = 3,17

$$\begin{aligned}\sigma^2 &= \frac{(2 - 3,17)^2 + (2 - 3,17)^2 + (3 - 3,17)^2 + (3 - 3,17)^2}{6 - 1} + \\ &+ \frac{(4 - 3,17)^2 + (5 - 3,17)^2}{5} = \\ &= \frac{1,3689 + 1,3689 + 0,0289 + 0,0289 + 0,6889 + 3,3489}{5} = \frac{6,8334}{5} = \\ &= 1,367\end{aligned}$$

Variabilita dat – směrodatná odchylka

směrodatná odchylka

{2,2,3,3,4,5}

průměr = 3,17

$$\sigma = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2} = 1,169$$

Variabilita dat – směrodatná odchylka

{2,2,3,3,4,5}	průměr = 3,17	SD = 1,17
{2,2,3,3,4,20}	průměr = 5,67	SD = 7,06
{5,5,6,6,6,6}	průměr = 5,67	SD = 0,52

Variabilita dat – směrodatná odchylka

{2,2,3,3,4,5}

průměr = 3,17

SD = 1,17

{2,2,3,3,4,20}

průměr = 5,67

SD = 7,06

{5,5,6,6,6,6}

průměr = 5,67

SD = 0,52

Modelový příklad

- průměrné relativní frekvence

$$\bar{f}(a) = 2, \bar{f}(na) = 2$$

- směrodatné odchylky

$$\sigma(a) = 1, \sigma(na) = 2$$

Modelový příklad

- z-skóre

$$z = \frac{\textit{rel. frekv. slova} - \textit{průměrná rel. frekv. slova ve vzorku}}{\textit{směrodatná odchylka slova}}$$

Modelový příklad

- z-skóre

$$z_{(\textit{text } A)}(a) = \frac{f_{\textit{rel}(\textit{text } A)}(a) - \bar{f}(a)}{\sigma(a)} = \frac{7 - 2}{1} = 5$$

Modelový příklad

- z-skóre

$$z_{(\textit{text A})}(a) = \frac{f_{\textit{rel}(\textit{text A})}(a) - \bar{f}(a)}{\sigma(a)} = \frac{7 - 2}{1} = 5$$

$$z_{(\textit{text B})}(a) = \frac{3 - 2}{1} = 1$$

$$z_{(\textit{text A})}(na) = \frac{4 - 2}{2} = 1$$

$$z_{(\textit{text B})}(na) = \frac{10 - 2}{1} = 4$$

Modelový příklad

- Delta vzdálenost mezi texty

$$\Delta_{(textA,textB)} = \frac{1}{N} \sum_{i=1}^N |z_i(textA) - z_i(textB)|$$

N ...počet analyzovaných slov i

Modelový příklad

- Delta vzdálenost mezi texty

$$\Delta_{(textA,textB)} = \frac{1}{N} \sum_{i=1}^N |z_i(textA) - z_i(textB)|$$

$$\begin{aligned} \Delta_{(textA,textB)} &= \frac{|z_a(textA) - z_a(textB)| + |z_{na}(textA) - z_{na}(textB)|}{2} = \\ &= \frac{|5 - 1| + |1 - 4|}{2} = \frac{7}{2} = 3.5 \end{aligned}$$

Vzdálenost mezi texty

$$\Delta_{(AB)} = \frac{1}{n} \sum_{i=1}^n \left| \frac{A_i - \mu_i}{\sigma_i} - \frac{B_i - \mu_i}{\sigma_i} \right|$$

n ... the number of MFW

A, B ... texts for the comparison

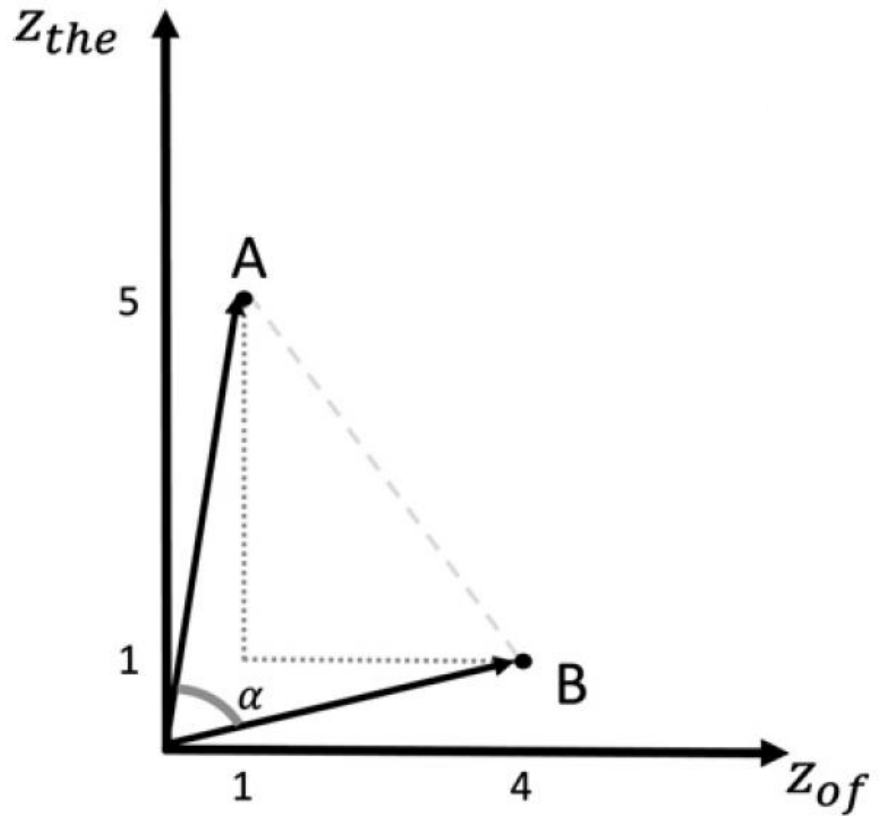
A_i ... the relative frequency of a given word in the text A

B_i ... the relative frequency of a given word in the text B

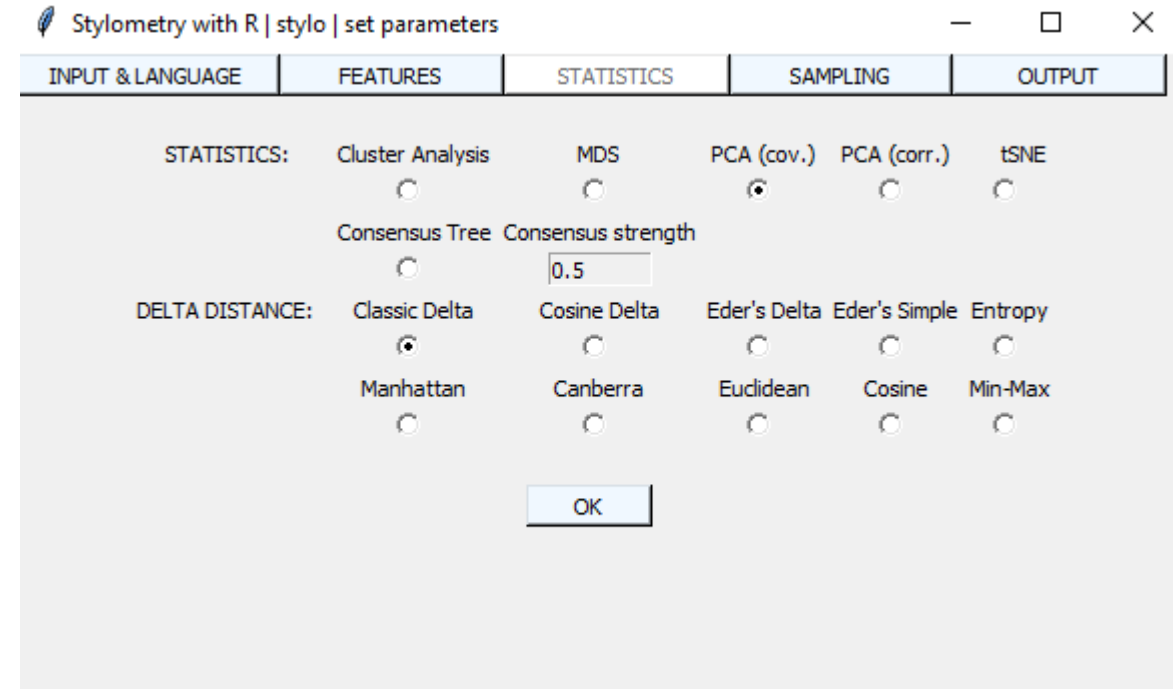
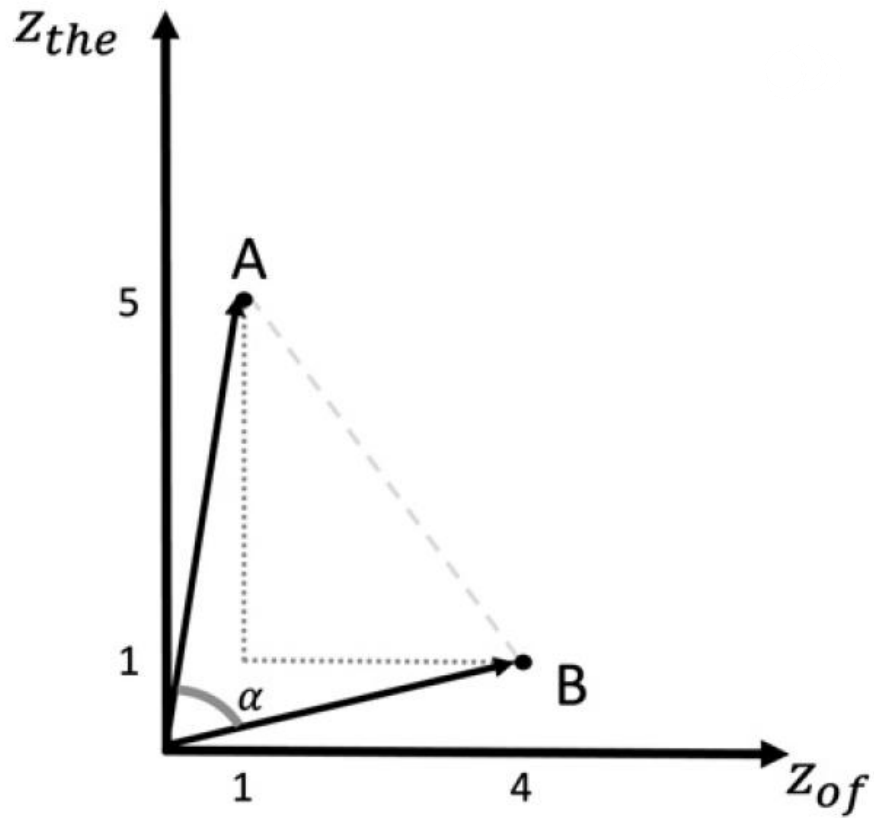
μ_i ... the average relative frequency of a given word in sample

σ_i ... the standard deviation of the relative frequency of a given word

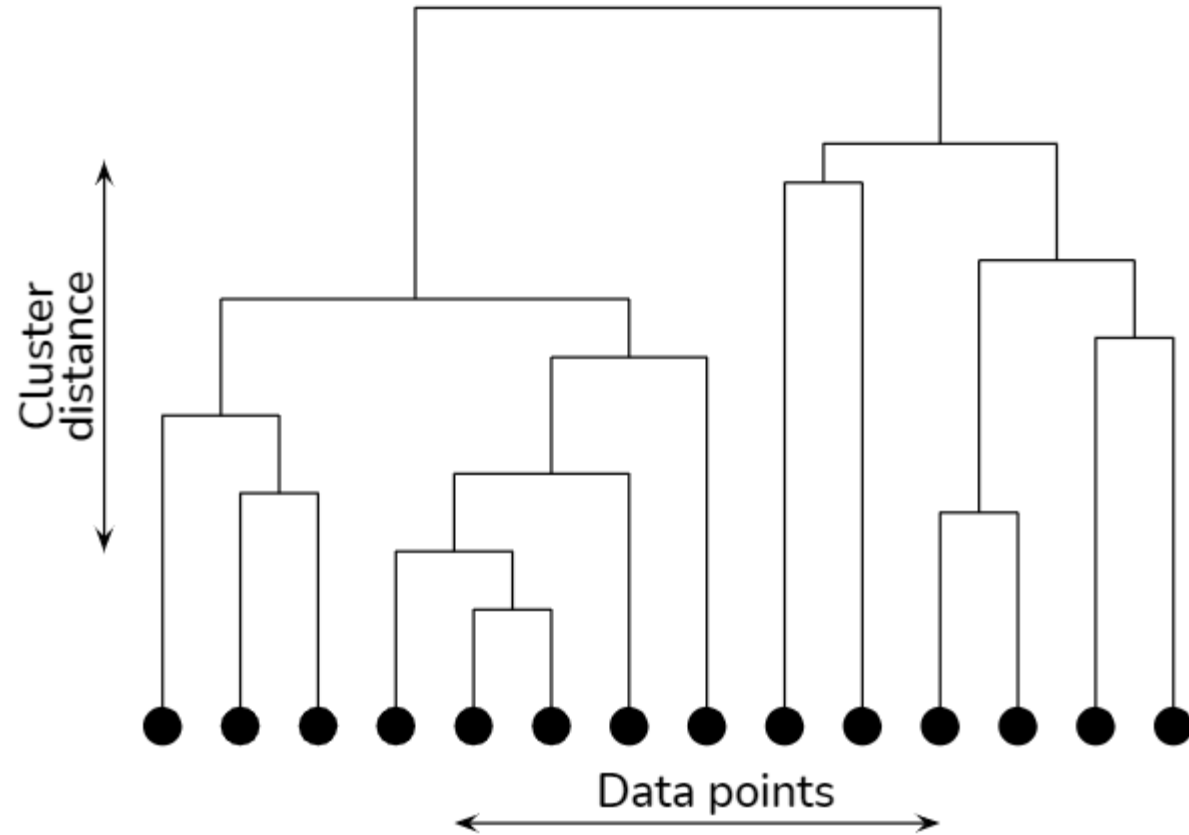
Vzdálenost mezi texty



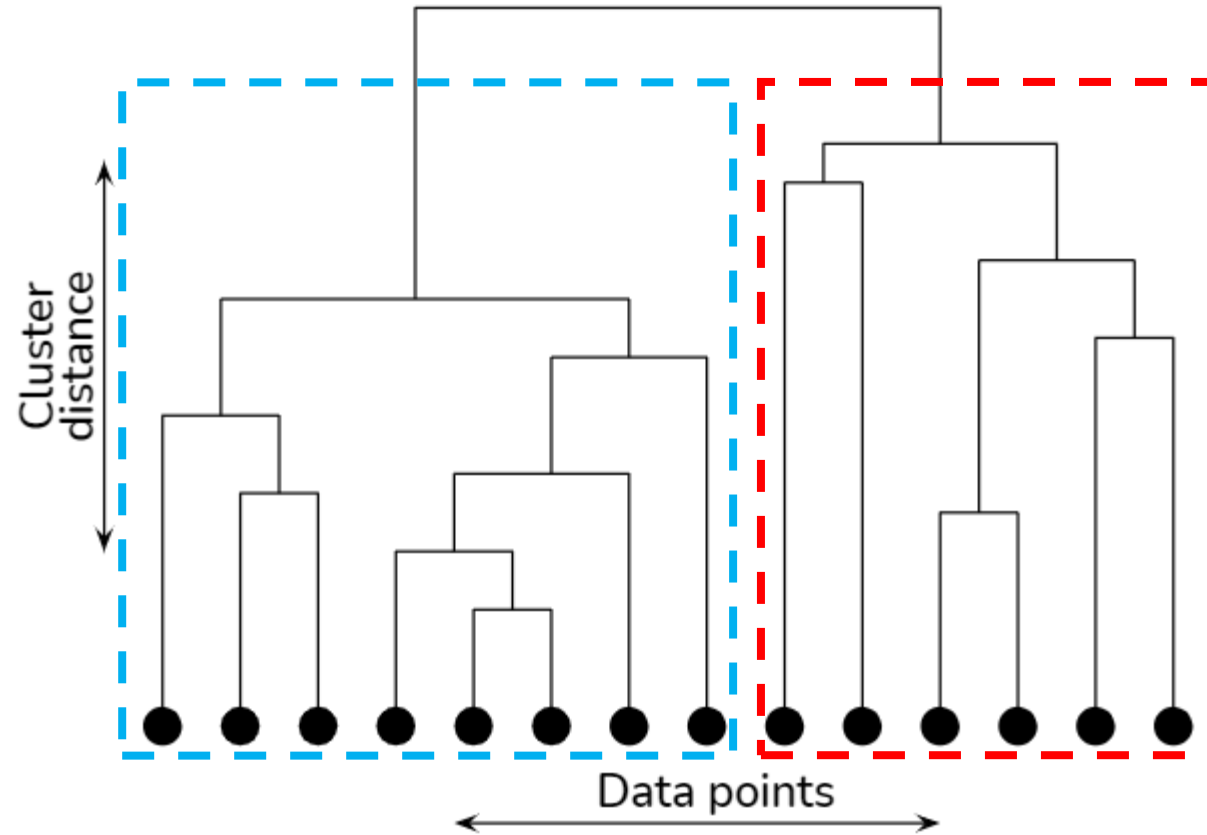
Vzdálenost mezi texty



Shluková analýza

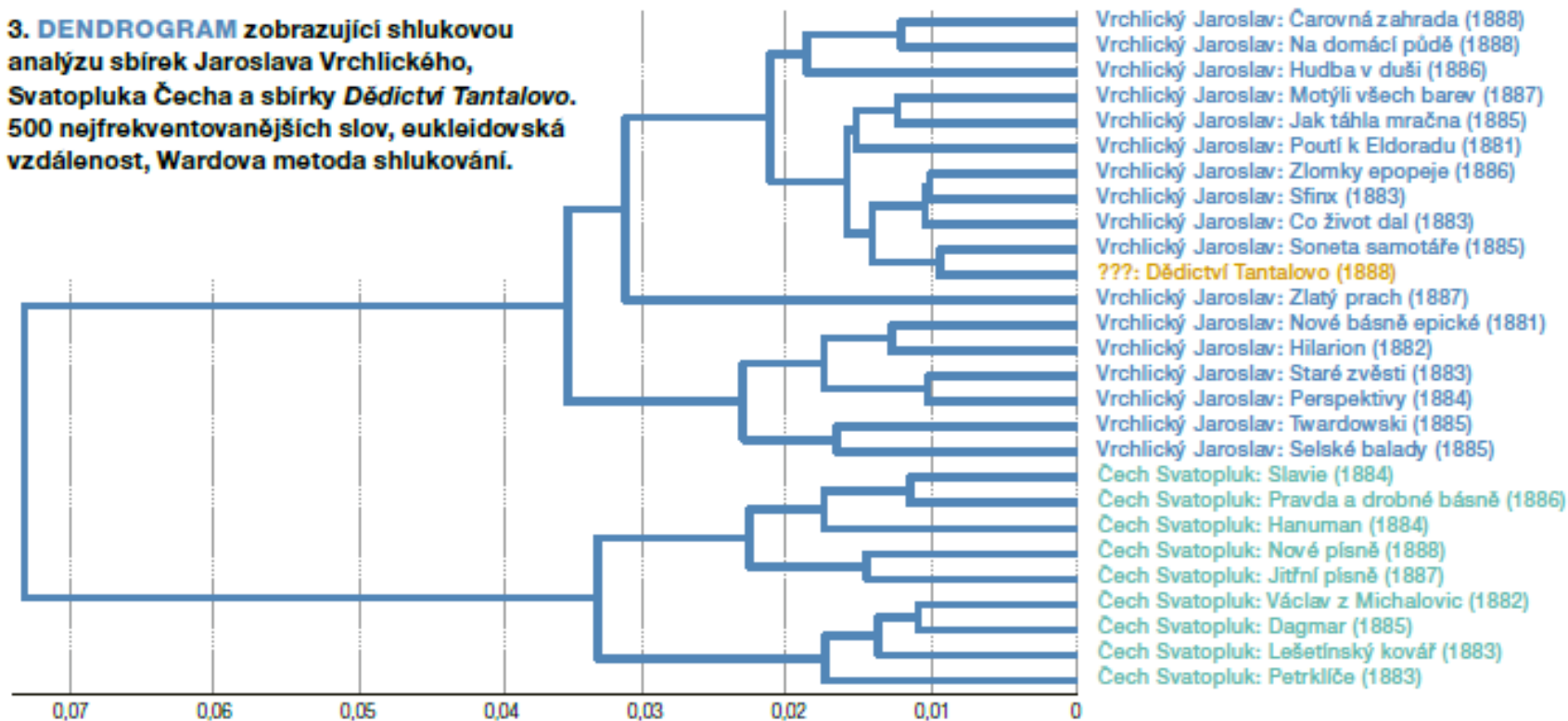


Shluková analýza



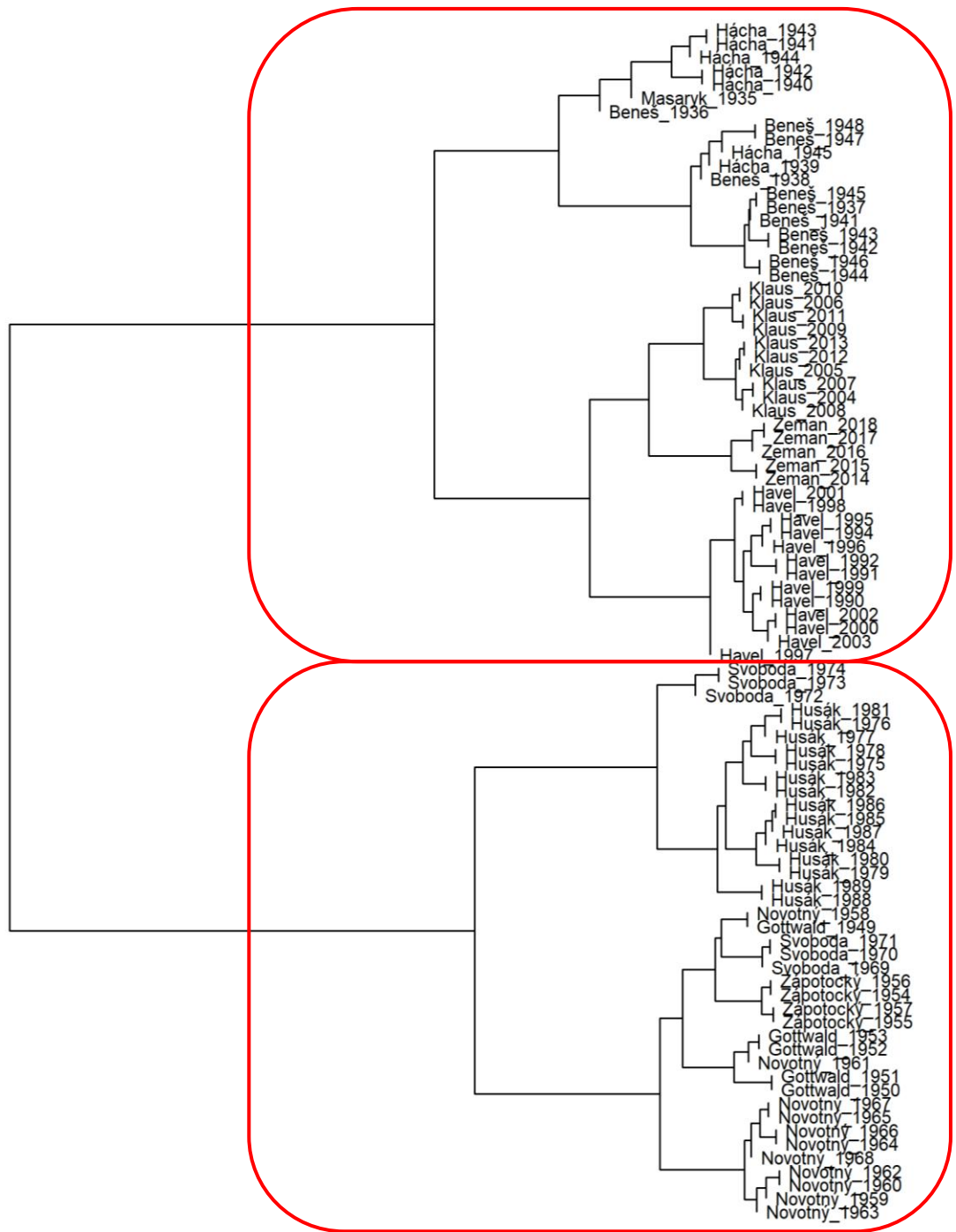
Shluková analýza

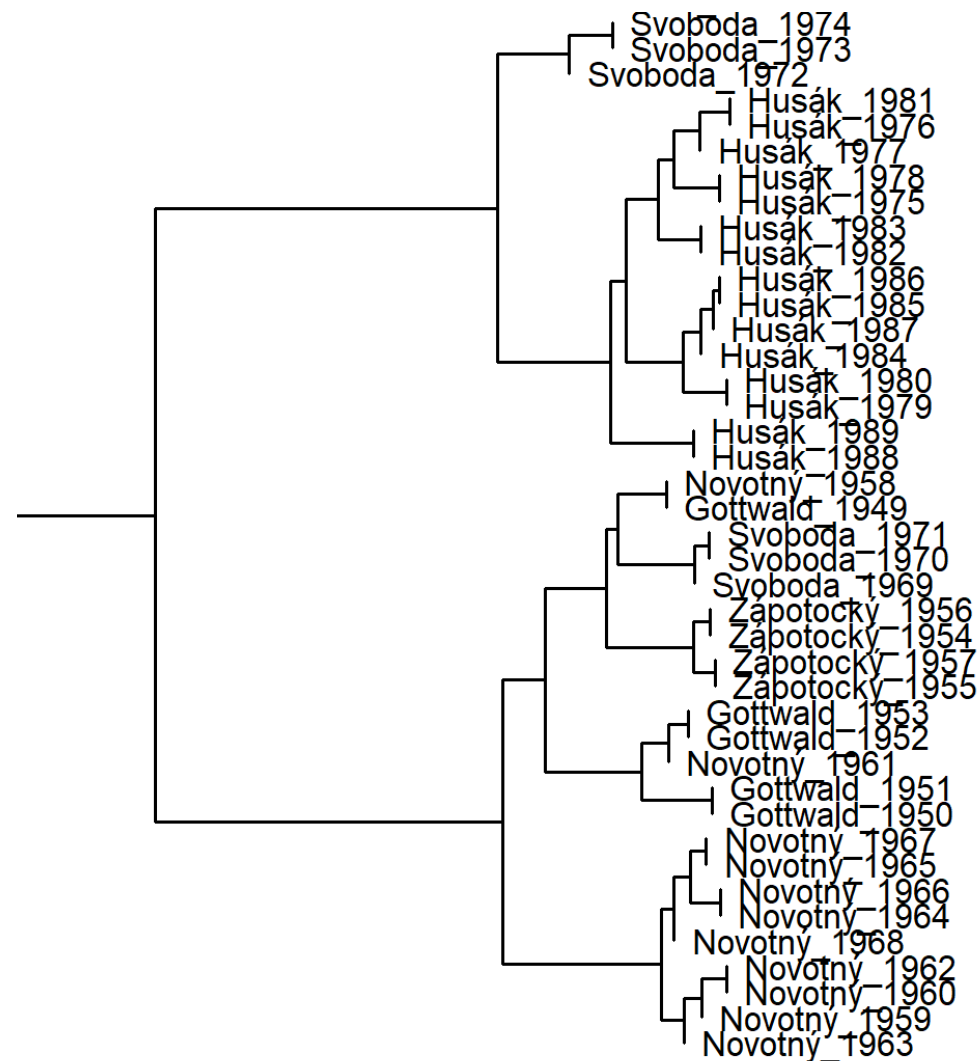
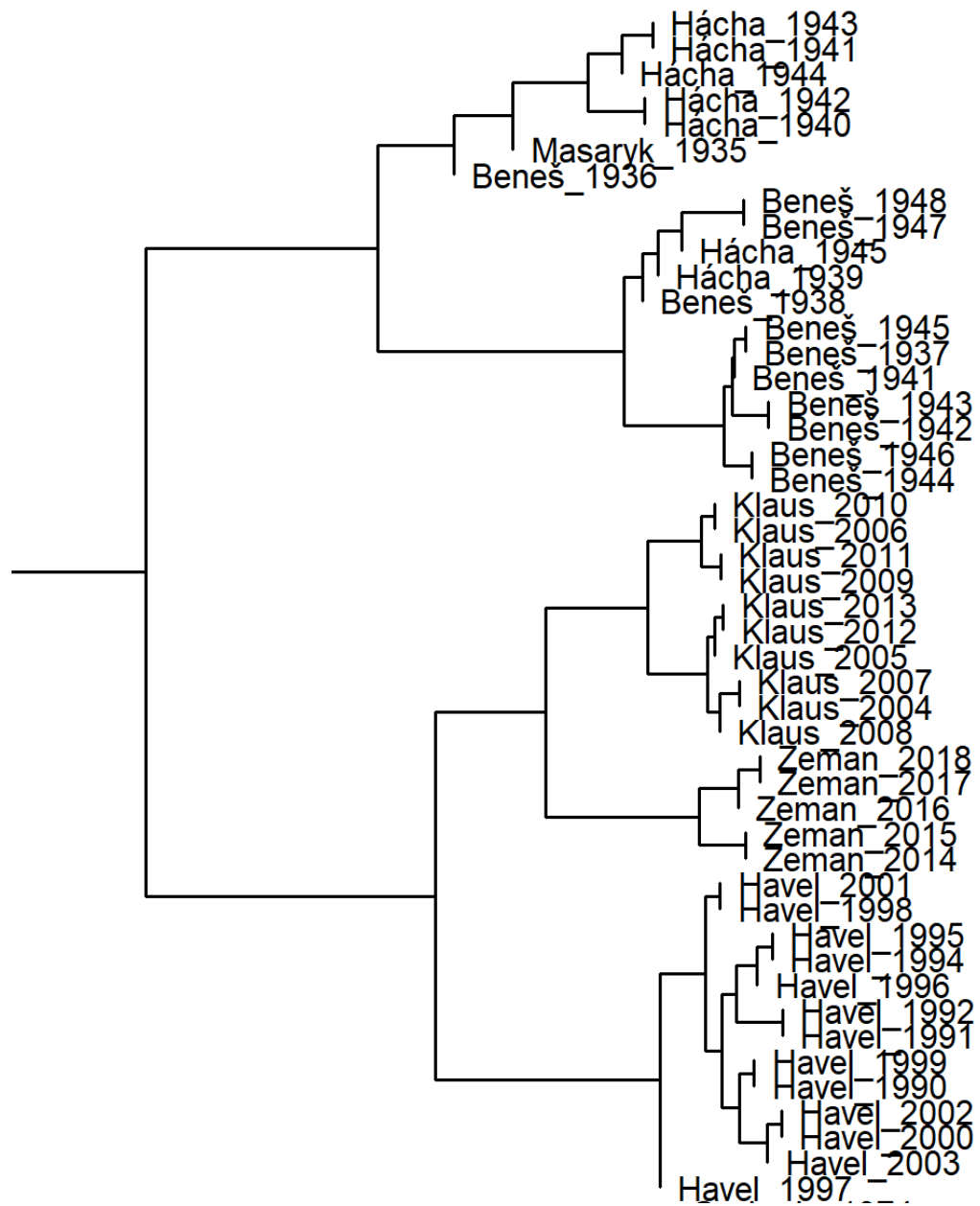
3. DENDROGRAM zobrazující shlukovou analýzu sbírek Jaroslava Vrchlického, Svatopluka Čecha a sbírky *Dědictví Tantalovo*. 500 nejfrekventovanějších slov, eukleidovská vzdálenost, Wardova metoda shlukování.

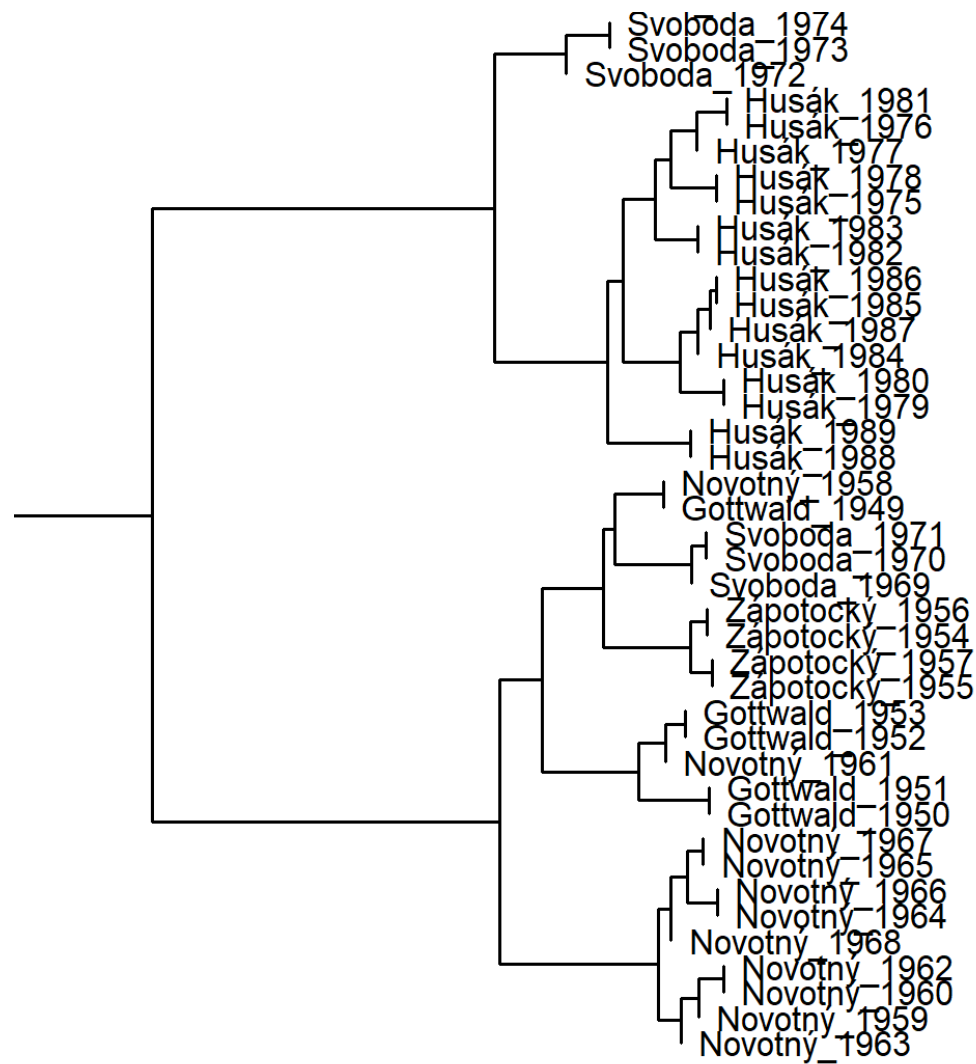
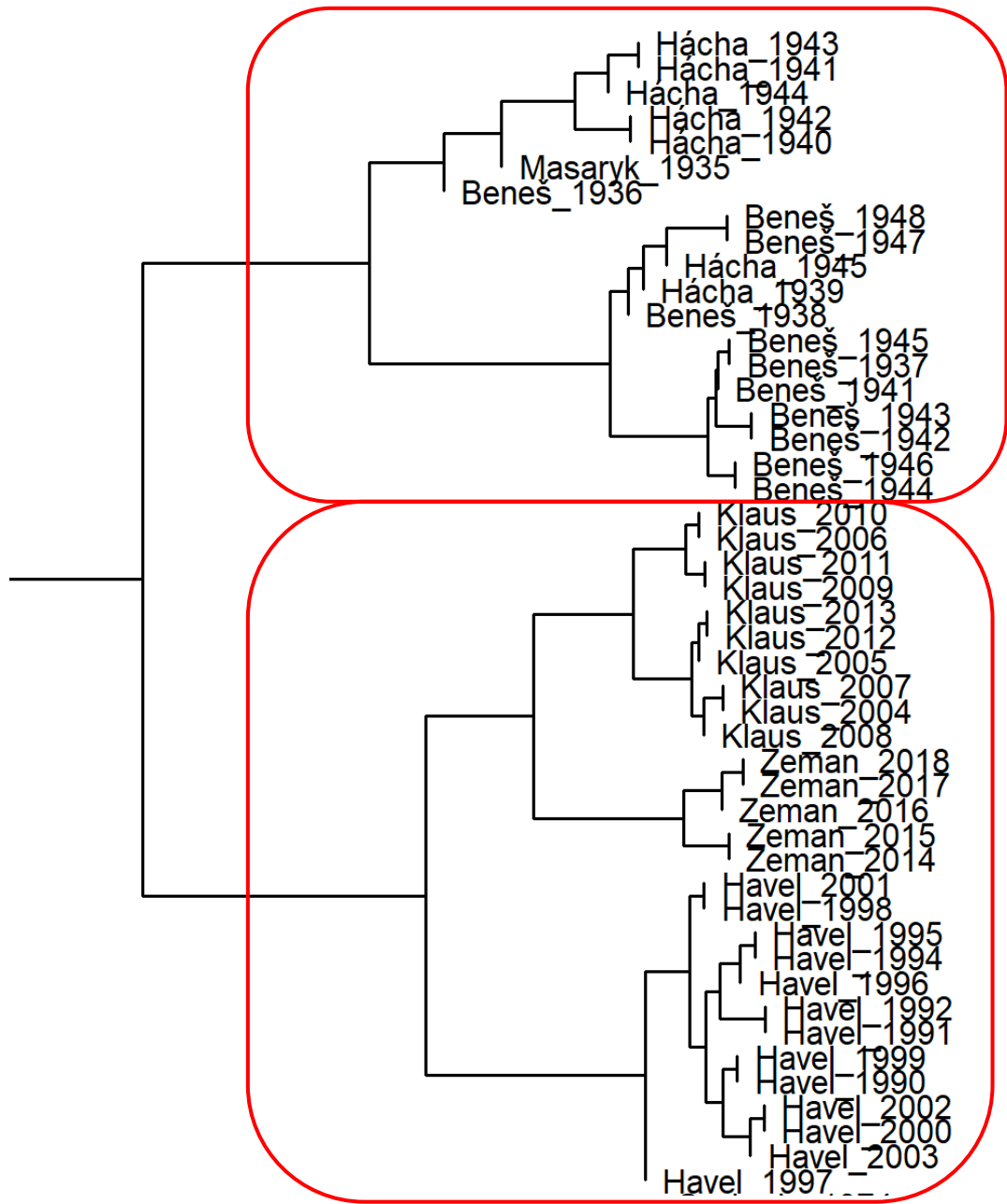


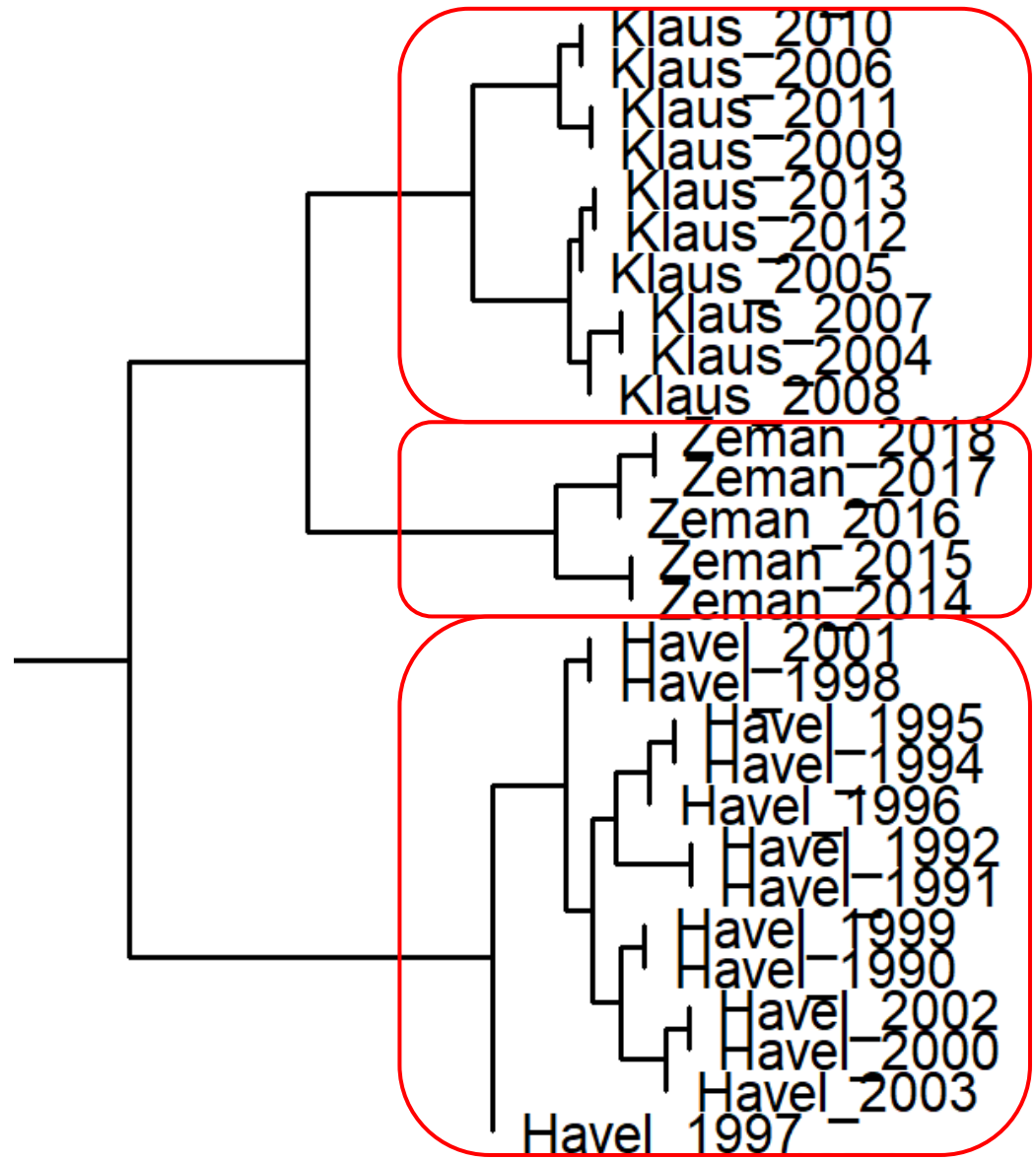
Prezidentské projevy

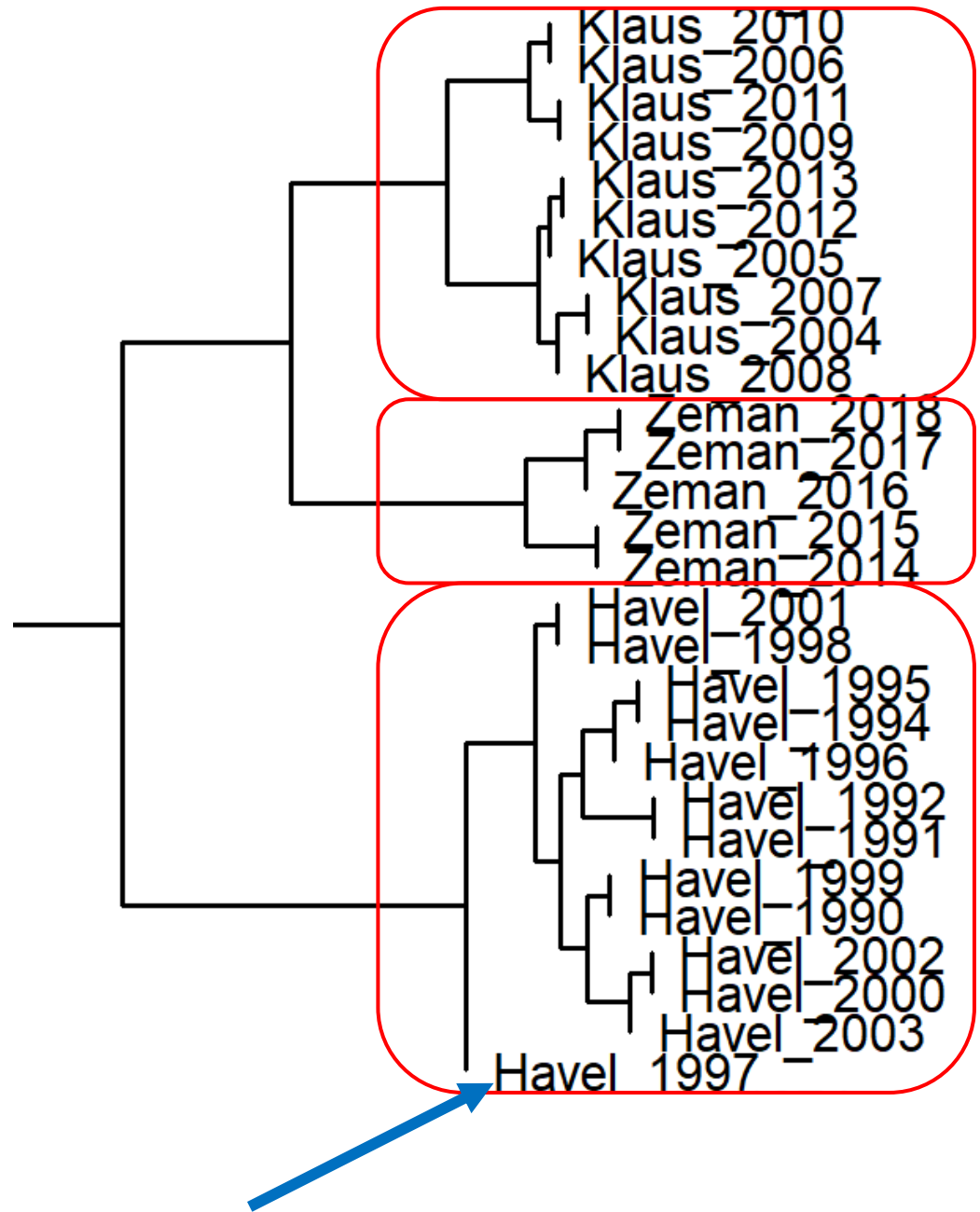
- Kubát, M., Mačutek, J., Čech, R. (2021). Communists spoke differently. An analysis of Czechoslovak and Czech annual presidential speeches. *Digital Scholarship in the Humanities*, 36, 138-152.
- presidential speeches: 1935–2018
- 100 MFW, culling = 60 %

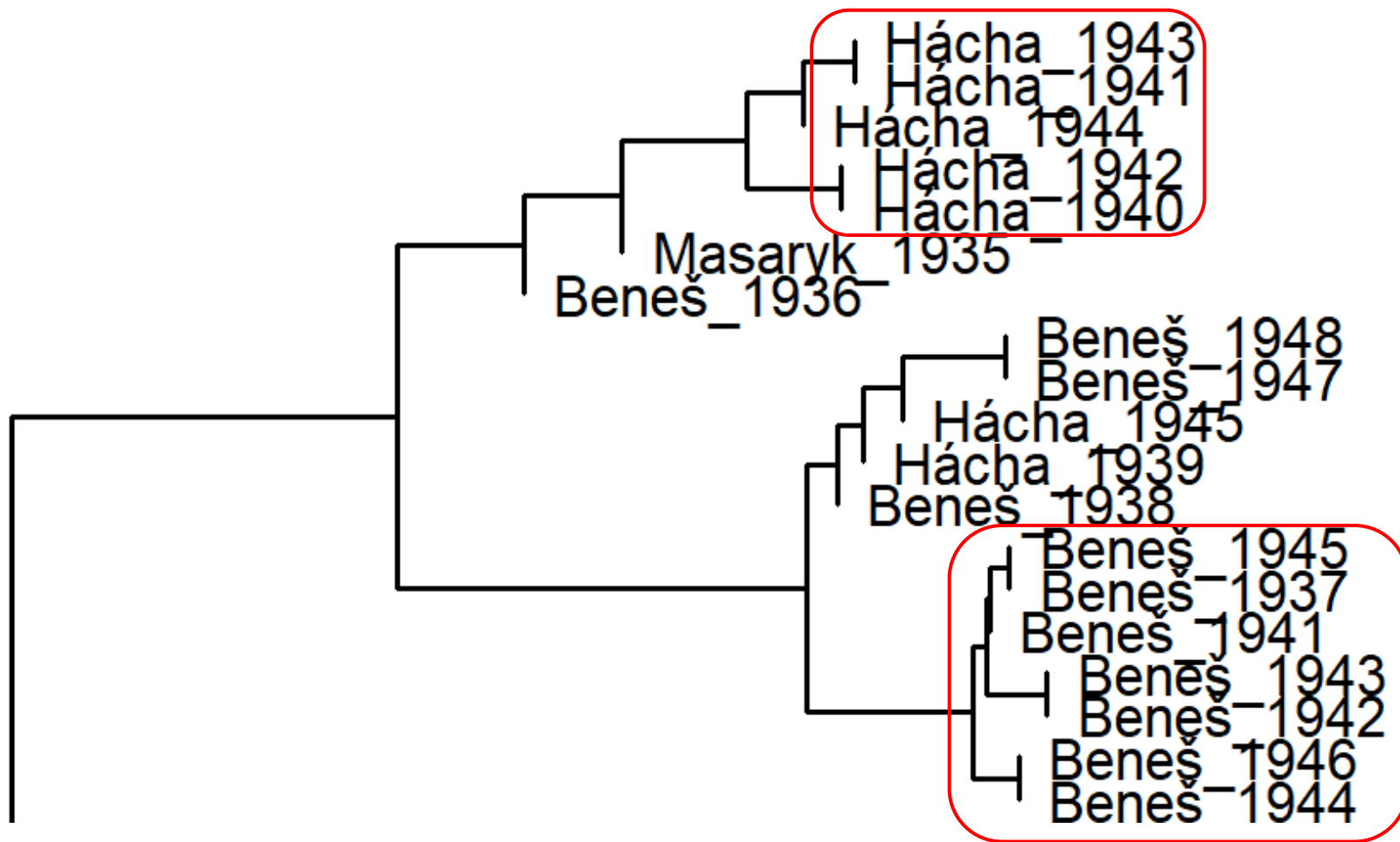


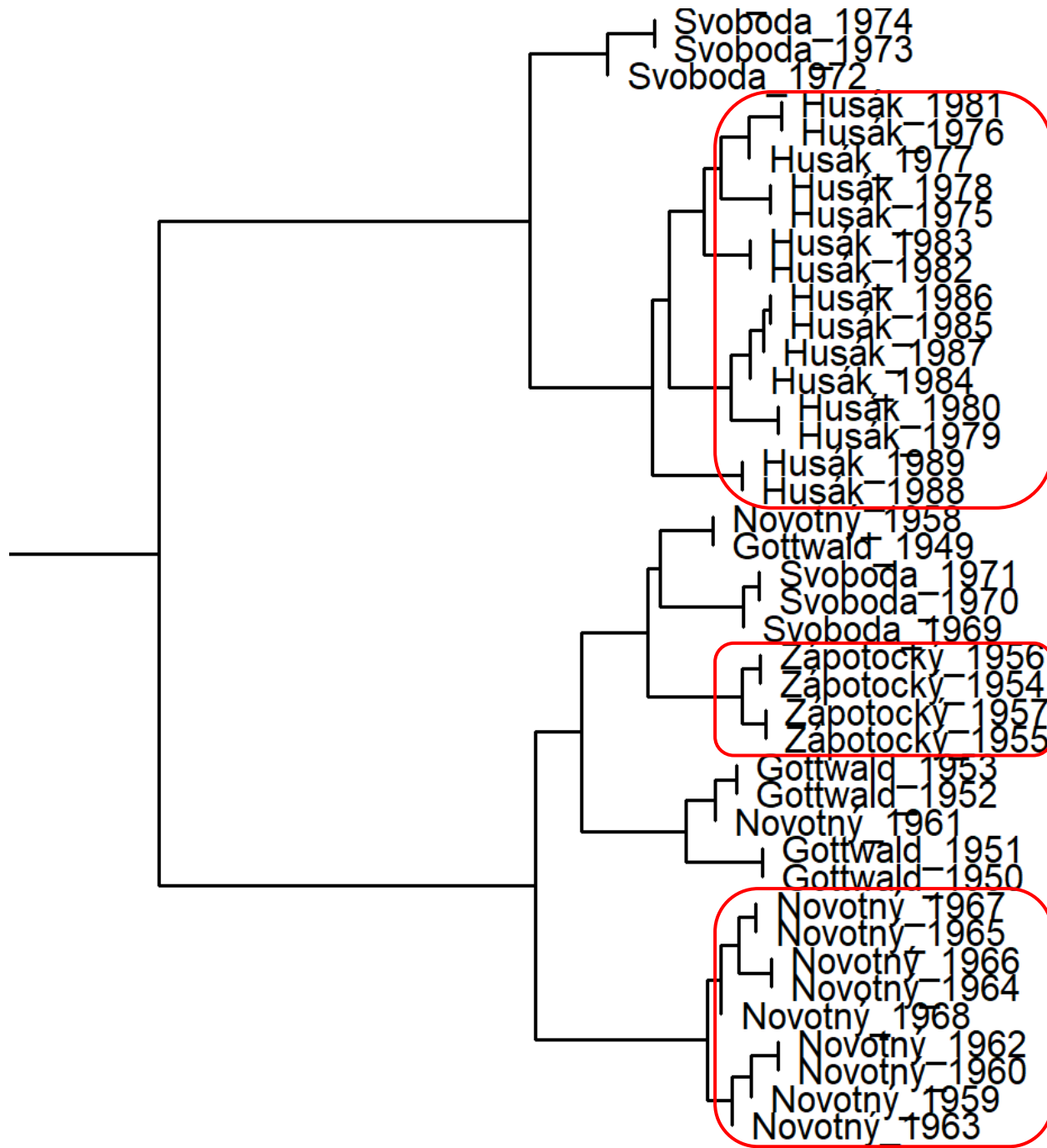


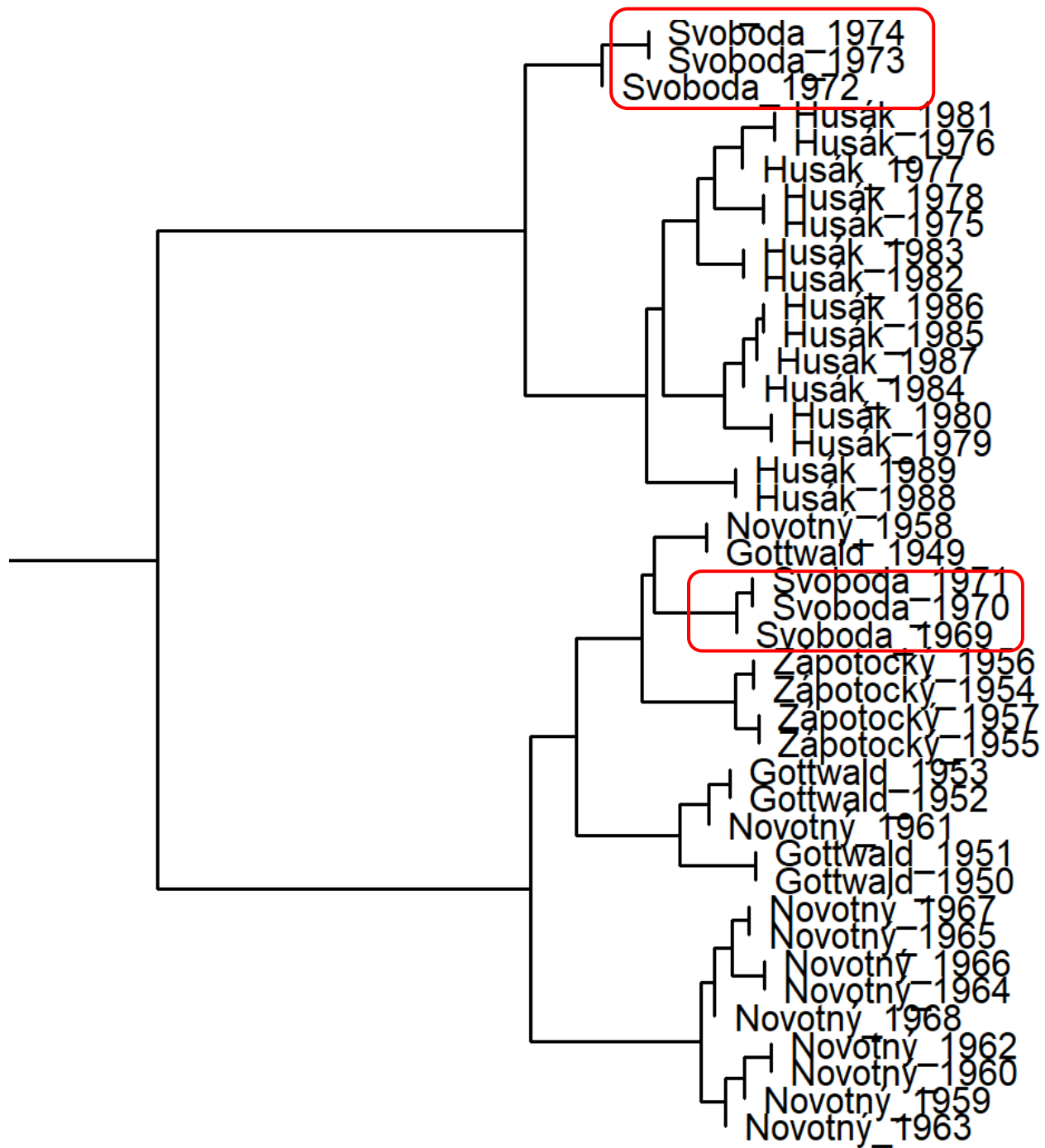












Bible svatováclavská & comments

- Kosek, P., Čech, R. (2018). Stylové aspekty Bible svatováclavské – stylometrická analýza. In Zand, G., Newerkla, S.M. (eds.). Jezuitská kultura v českých zemích / Jesuitische Kultur in den böhmischen Ländern. Host, 195-209.

Bible svatováclavská & comments

- New Testament
 - Konstanc → Šteyer
- Old Testament
 - Šteyer → Barner
 - Job – divide
- „[...] (*domníváme se, že by se stylistickým rozdílem textu dalo zjistit, odkud překládal už jen Šteyer*) [...]“ VINTR (1997)

Bible svatováclavská & comments

- comments
 - Šteyer: New Testament + Genesis
 - Barner: Old Testament

Bible svatováclavská & comments

- translation
- sacred text
- editing

Bible svatováclavská & comments

- New Testament

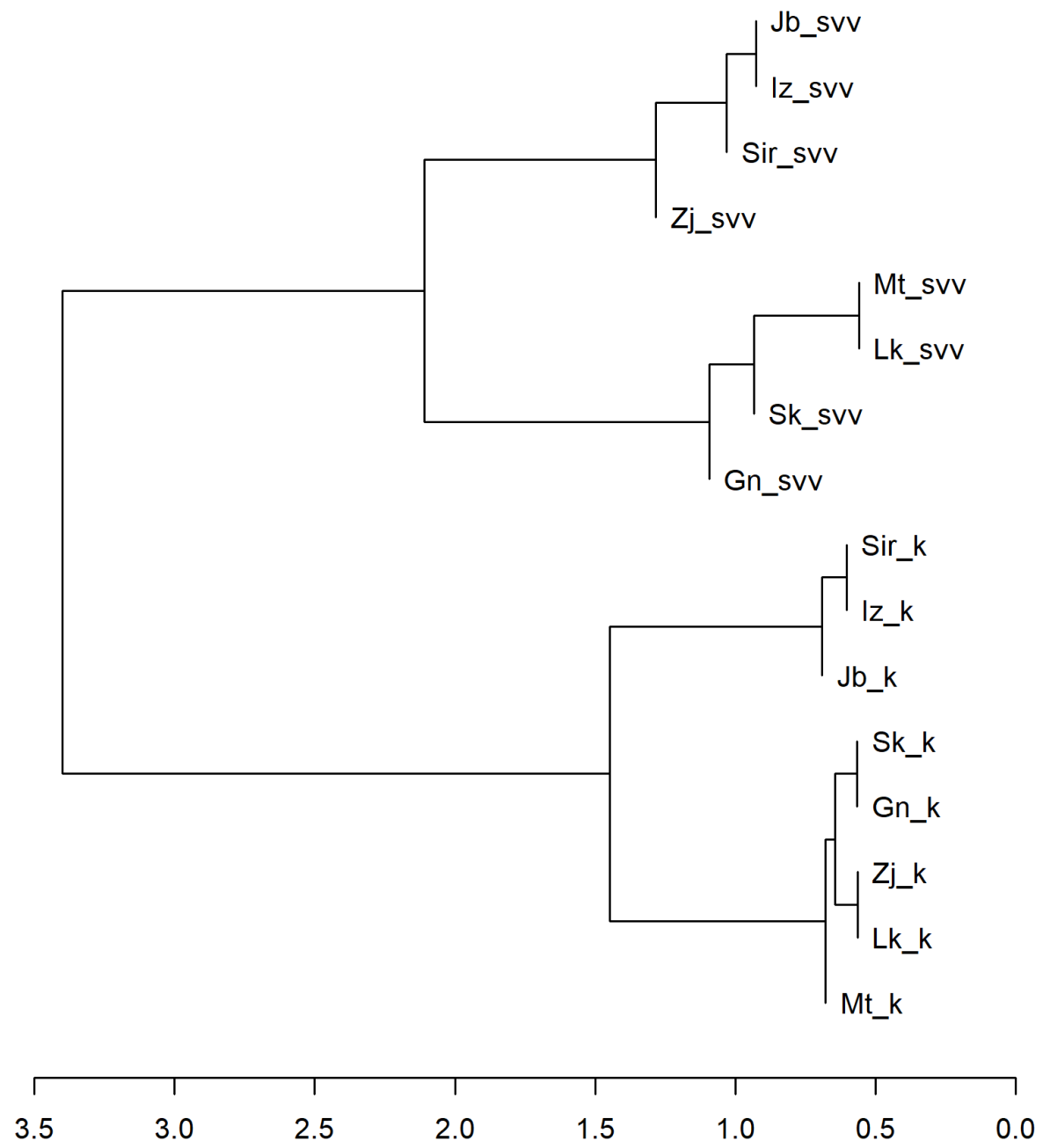
- Mt, Lk, Sk, Zj

- Old Testament

- Gn, Jb, Iz, Sir

- 100 MFW

- culling = 0



Jak to, že to funguje?

- Nini, s. 32

Stylo

- Eder, M., Rybicki, J., & Kestemont, M. (2016). Stylometry with R: a package for computational text analysis. *The R Journal*, 8(1).

Stylometry with R | stylo | set parameters

INPUT & LANGUAGE	FEATURES	STATISTICS	SAMPLING	OUTPUT	
INPUT:	plain text	xml	xml (plays)	xml (no titles)	html
	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
LANGUAGE:	English	English (contr.)	English (ALL)	Latin	Latin (u/v > u)
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	Polish	Hungarian	French	Italian	Spanish
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	Dutch	German	CJK	Other	Native encoding
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="checkbox"/>
<input type="button" value="OK"/>					

Stylo

Stylometry with R | stylo | set parameters

INPUT & LANGUAGE	FEATURES	STATISTICS	SAMPLING	OUTPUT
FEATURES:	words <input checked="" type="radio"/>	chars <input type="radio"/>	ngram size <input type="text" value="1"/>	preserve case <input type="checkbox"/>
MFV SETTINGS:	Minimum <input type="text" value="100"/>	Maximum <input type="text" value="100"/>	Increment <input type="text" value="100"/>	Start at freq. rank <input type="text" value="1"/>
CULLING:	Minimum <input type="text" value="0"/>	Maximum <input type="text" value="0"/>	Increment <input type="text" value="20"/>	List Cutoff <input type="text" value="5000"/>
				Delete pronouns <input type="checkbox"/>
VARIOUS:	Existing frequencies <input type="checkbox"/>	Existing wordlist <input type="checkbox"/>	Select files manually <input type="checkbox"/>	List of files <input type="checkbox"/>

OK

Stylometry with R | stylo | set parameters

INPUT & LANGUAGE	FEATURES	STATISTICS	SAMPLING	OUTPUT	
STATISTICS:	Cluster Analysis <input type="radio"/>	MDS <input type="radio"/>	PCA (cov.) <input checked="" type="radio"/>	PCA (corr.) <input type="radio"/>	tSNE <input type="radio"/>
	Consensus Tree <input type="radio"/>	Consensus strength <input type="text" value="0.5"/>			
DELTA DISTANCE:	Classic Delta <input checked="" type="radio"/>	Cosine Delta <input type="radio"/>	Eder's Delta <input type="radio"/>	Eder's Simple <input type="radio"/>	Entropy <input type="radio"/>
	Manhattan <input type="radio"/>	Canberra <input type="radio"/>	Euclidean <input type="radio"/>	Cosine <input type="radio"/>	Min-Max <input type="radio"/>

OK