

Analýza klíčových slov versus tematická koncentrace textu

Radek Čech
(KČJ FF OU)

ÚČNK FF UK
Praha
3. 12. 2013

Analýza klíčových slov

- cíl: detekovat slova, která lze díky relativně vysoké/nízké frekvenci považovat za slova vyjadřující hlavní téma/témata textu
- klíčové slovo – slovo, které se v daném textu používá se signifikantně vyšší nebo nižší frekvencí, než je tomu v referenčním korpusu
- testovým kritériem
 - chí-kvadrát test
 - log-likelihood test

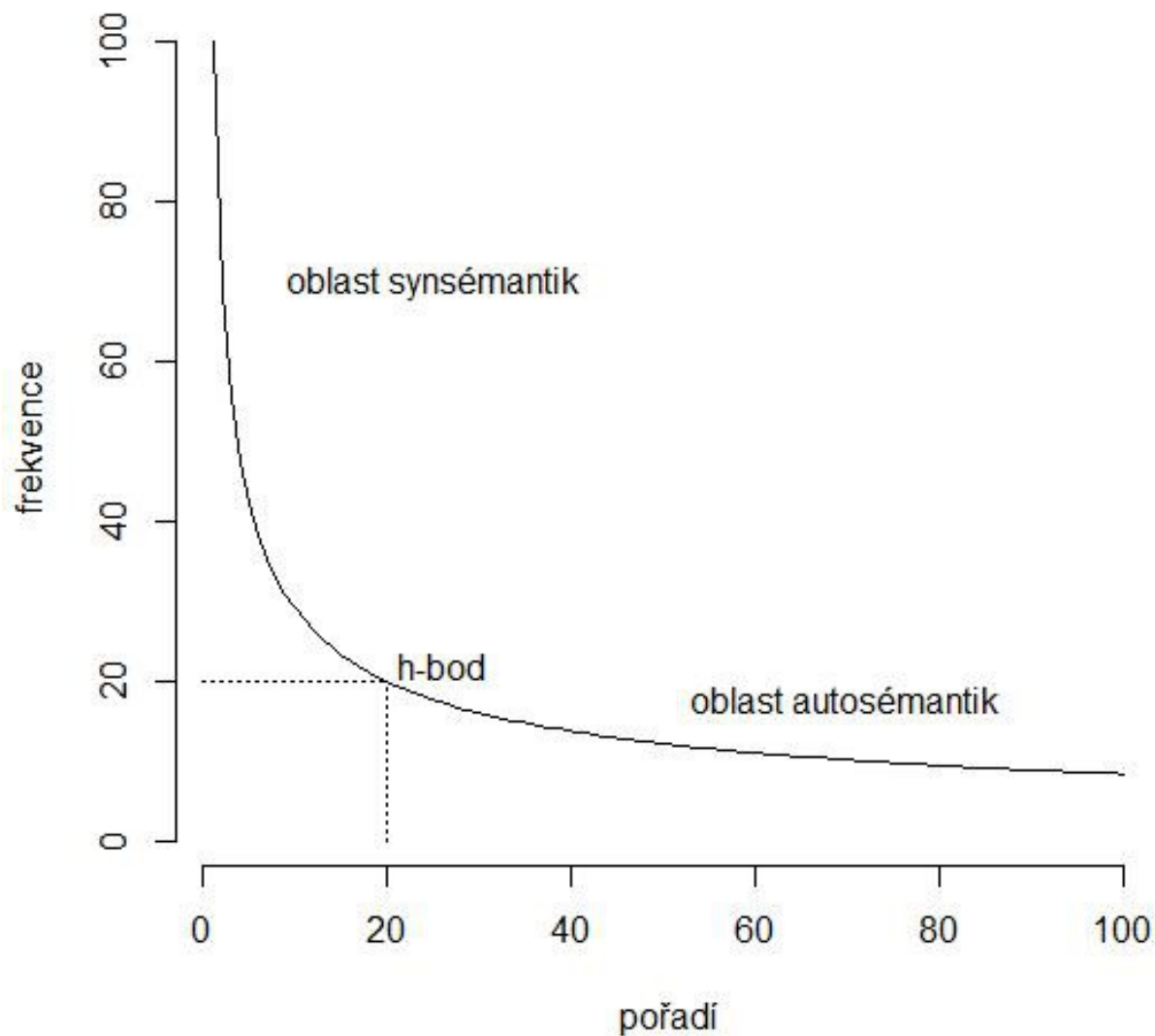
- L. Jehlička *Mládí a poezie* (log-likelihood; $p = 0,001$; $\min(f) = 3$; all signif. types; stop-list: pron., prep, conj., number; ignore case; ref. corpus: SYN2010)

Form	LL	Fq(text)	Fq(ref)	Form	LL	Fq(text)	Fq(ref)
<i>hrdinstvím</i>	85.144	5	38	<i>jest</i>	31.067	4	3346
<i>hrdinství</i>	77.212	6	386	<i>poezii</i>	27.974	3	1146
<i>mládí</i>	60.701	7	3713	<i>ducha</i>	21.622	3	3328
<i>poezie</i>	46.29	5	1975	<i>mladí</i>	19.363	3	4873
<i>katolictví</i>	41.943	3	110	<i>člověku</i>	18.805	3	5356
<i>otroctví</i>	33.223	3	476	<i>lásky</i>	18.012	3	6128
<i>nejen</i>	31.731	7	30472	<i>především</i>	17.619	5	37073
<i>oběti</i>	31.342	4	3232	<i>pak</i>	17.489	9	161980
<i>životu</i>	31.223	4	3281	<i>není</i>	16.93	8	129953

Tematická koncentrace textu

- cíl: kvantifikovat celkovou charakteristiku textu, přičemž kvantifikace se chápe jako prostředek, který umožní testovat rozdíly mezi jednotlivými texty (příp. skupinami textů); tyto rozdíly mohou být důsledkem vliv žánrů, autorství atd.
- detekce jednotlivých slov a kvantifikace jejich tzv. tematické váhy „jen“ dílčím krokem
- TK založena na frekvenční struktuře textu a tzv. h-bodu

Tematická koncentrace textu



- L. Jehlička *Mládí a poezie*; h -bod = 9; sekundární TK: $2h = 18$

pořadí	slovo	frekvence	pořadí	slovo	frekvence
1	<i>a</i>	53	11	<i>na</i>	8
2	<i>v</i>	21	12	<i>li</i>	8
3	<i>je</i>	19	13	<i>není</i>	8
4	<i>k</i>	15	14	<i>však</i>	8
5	<i>se</i>	15	15	<i>nejen</i>	7
6	<i>ale</i>	15	16	<i>mládí</i>	7
7	<i>že</i>	13	17	<i>s</i>	7
8	<i>aby</i>	10	18	<i>hrdinství</i>	6
9	<i>pak</i>	9	19	<i>bez</i>	6
10	<i>i</i>	9	20	<i>bylo</i>	6

- L. Jehlička *Mládí a poezie*; h -bod = 10,5

pořadí	lemma	frekvence	pořadí	lemma	frekvence
1	<i>být</i>	55	11	<i>ten</i>	10
2	<i>a</i>	53	12	<i>jenž</i>	10
3	<i>v</i>	24	13	<i>tento</i>	10
4	<i>se</i>	21	14	<i>aby</i>	10
5	<i>k</i>	15	15	<i>pak</i>	9
6	<i>ale</i>	15	16	<i>i</i>	9
7	<i>že</i>	13	17	<i>poezie</i>	8
8	<i>on</i>	13	18	<i>však</i>	8
9	člověk	13	19	<i>na</i>	8
10	hrdinství	11	20	<i>život</i>	8

Tematická váha slova

- tematická váha slova

$$(h - r') \cdot f(r')$$

- normalizace

$$C = \sum_{r=1}^h (h - r) \cdot f(1)$$

- index tematické váhy slova

$$TV_{slovo} = 2 \frac{(h - r') \cdot f(r')}{h(h - 1) \cdot f(1)}$$

Tematická váha slova

- index tematické váhy lemmatu *člověk*
 - $h = 10,5$; $r = 8$; $f = 13$

$$TV_{\text{člověk}} = 2 \frac{(h - r') \cdot f(r')}{h(h - 1) \cdot f(1)} = 2 \frac{(10,5 - 8) \cdot 13}{10,5(10,5 - 1) \cdot 55} = 0,011848$$

KW vs. TK

- společné znaky
 - analýza lexikálních charakteristik textu
 - intersubjektivita
 - automatická detekce jednotek reprezentujících hlavní téma/témata textu

Problémy KW & TK

- KW & TK
 - volba jazykových jednotek
 - slovní forma, lemma, koreferenční jednotka, hřeb...
 - vliv délky textu
- KW
 - vliv referenčního korpusu
 - volba hladiny významnosti
- TK
 - v mnohých případech se nevyskytuje nad h -bodem žádné autosémantické slovo

TK

- cíl: kvantifikovat celkovou charakteristiku textu, přičemž kvantifikace se chápe jako prostředek, který umožní testovat rozdíly mezi jednotlivými texty (příp. skupinami textů); tyto rozdíly mohou být důsledkem vliv žánrů, autorství atd.

$$TK_{text} = \sum_{r'=1}^T 2 \frac{(h-r') \cdot f(r')}{h(h-1) \cdot f(1)}$$

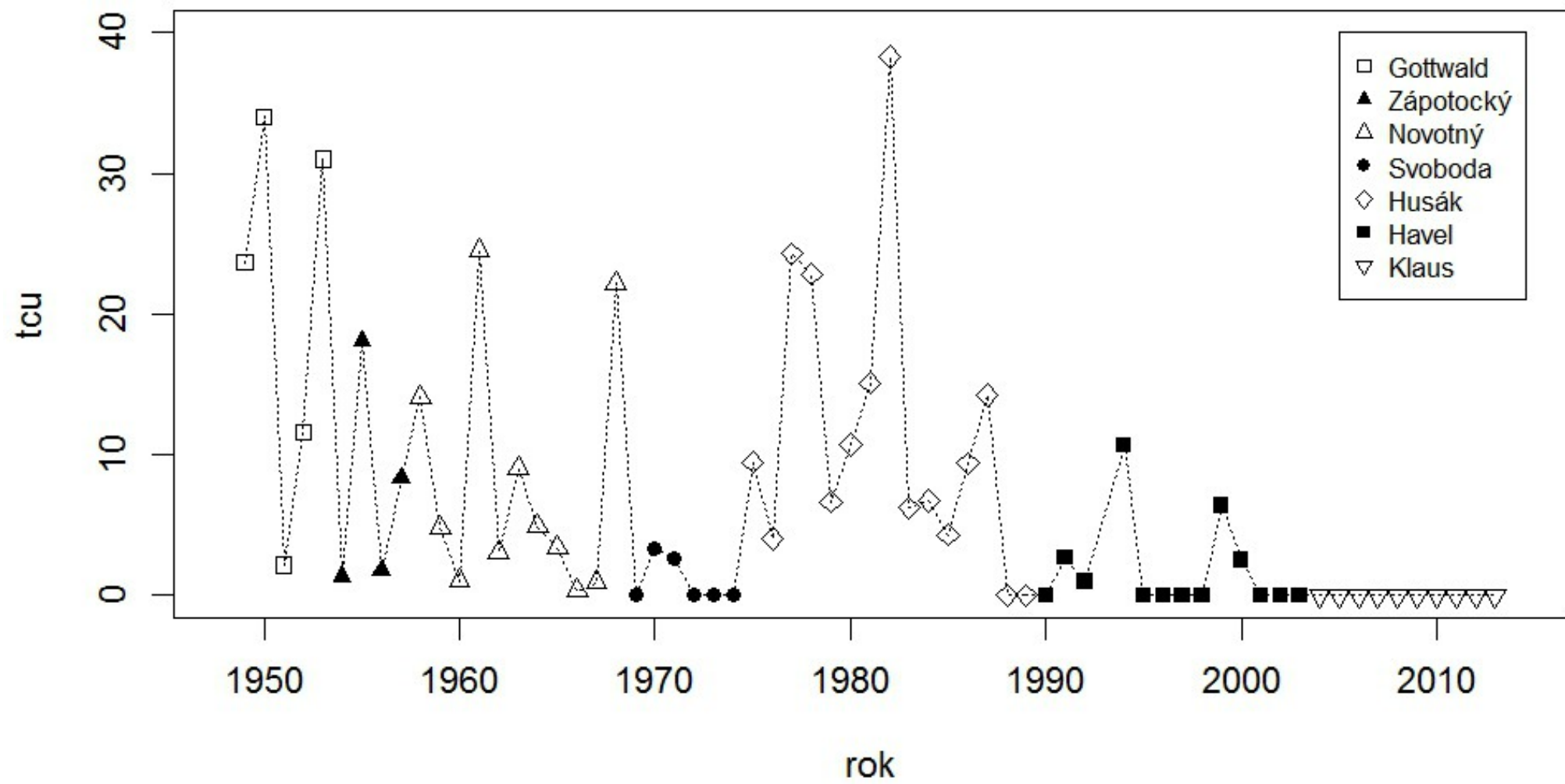
TK

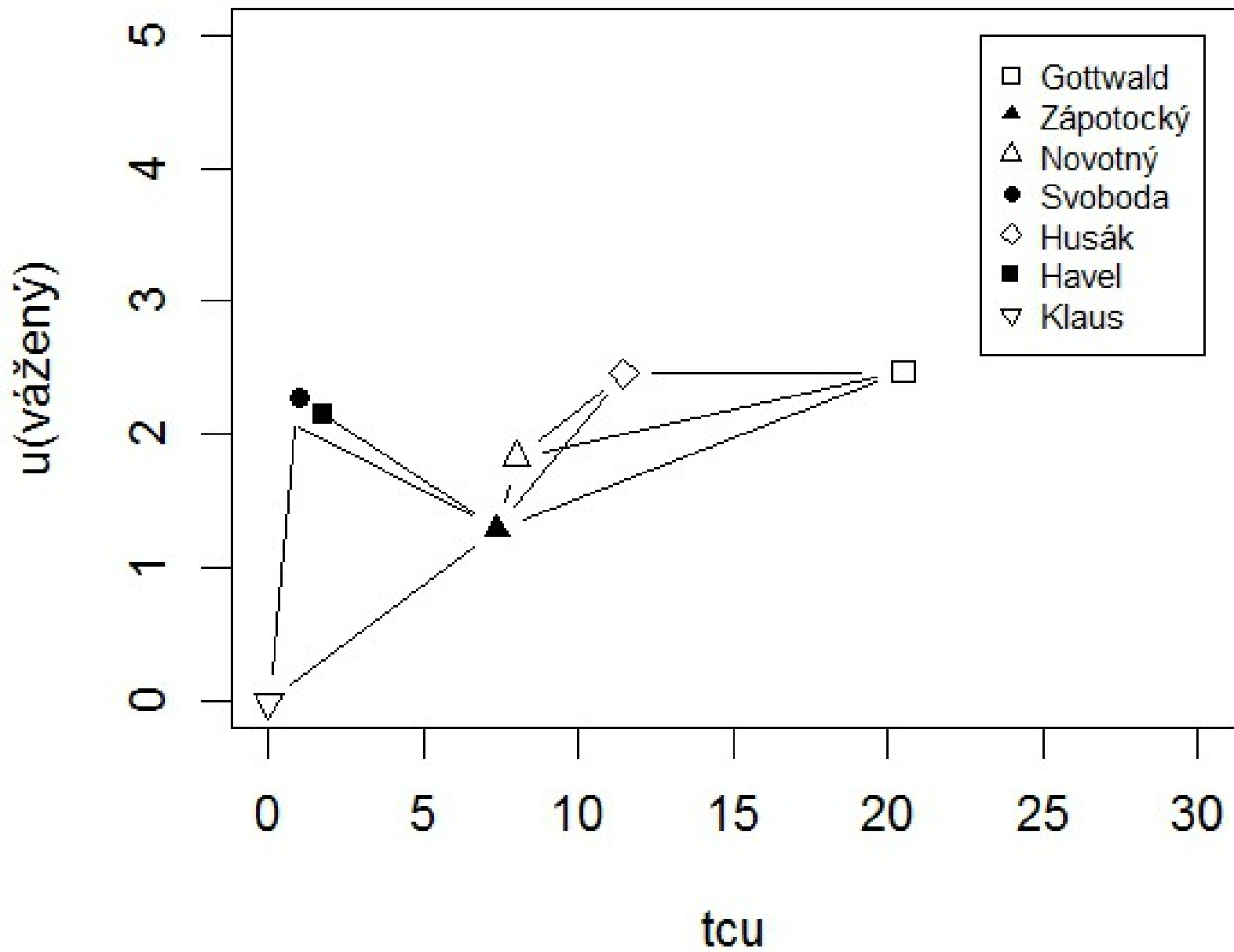
$$\text{VAR}(TK) = \left(\frac{2}{h(h-1)f(1)} \right)^2 \cdot \left(\sum_{r'=1}^T f(r') \right) \cdot m_{2,r'}$$

$$m_{2,r'} = \frac{\sum_{r'=1}^T (r' - m_{1,r'})^2 f(r')}{\sum_{r'=1}^T f(r')}$$

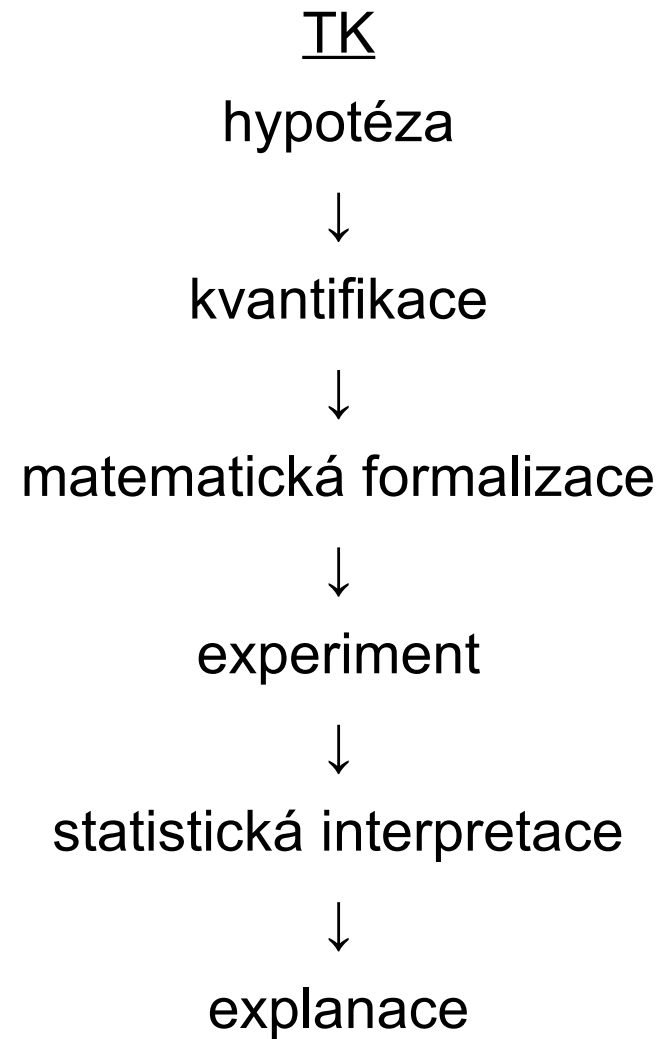
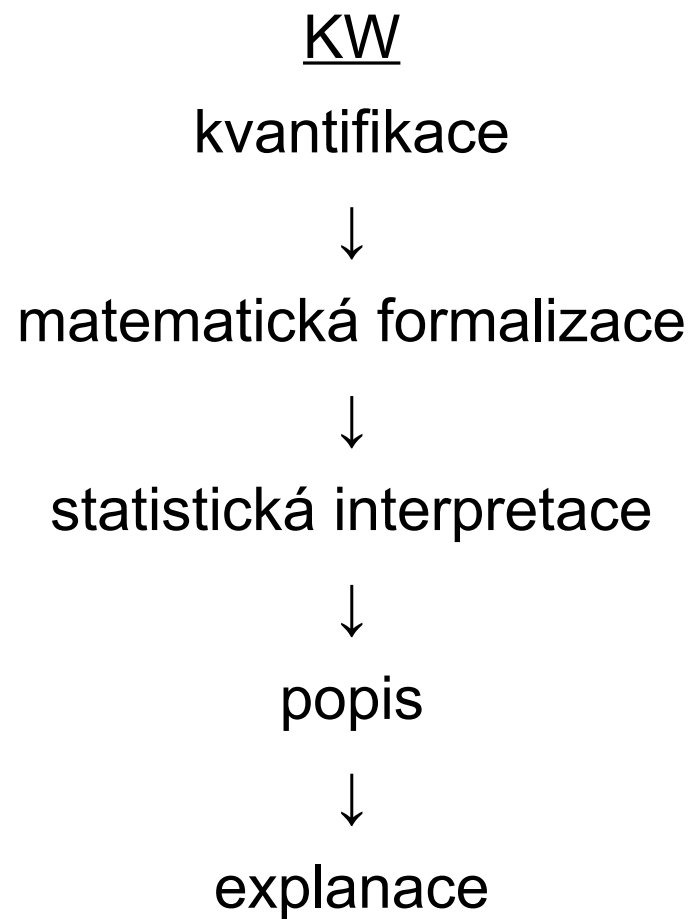
$$m_{1,r'} = \frac{\sum r' \cdot f(r')}{f(r')}$$

$$|u| = \frac{TK_1 - TK_2}{\sqrt{\text{VAR}(TK_1) + \text{VAR}(TK_2)}}$$

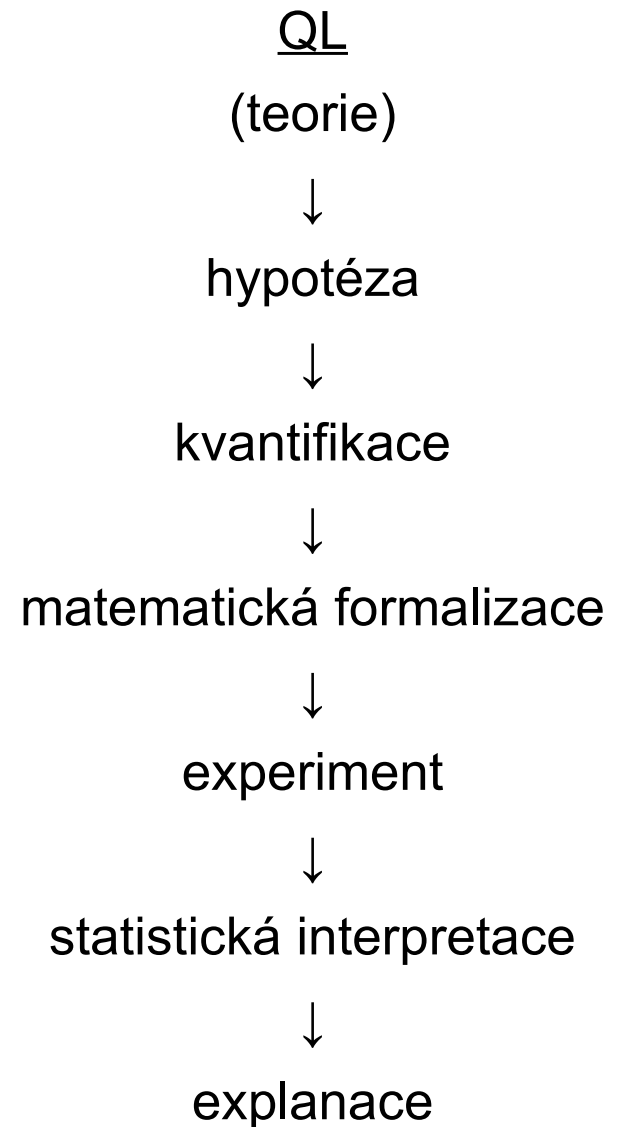
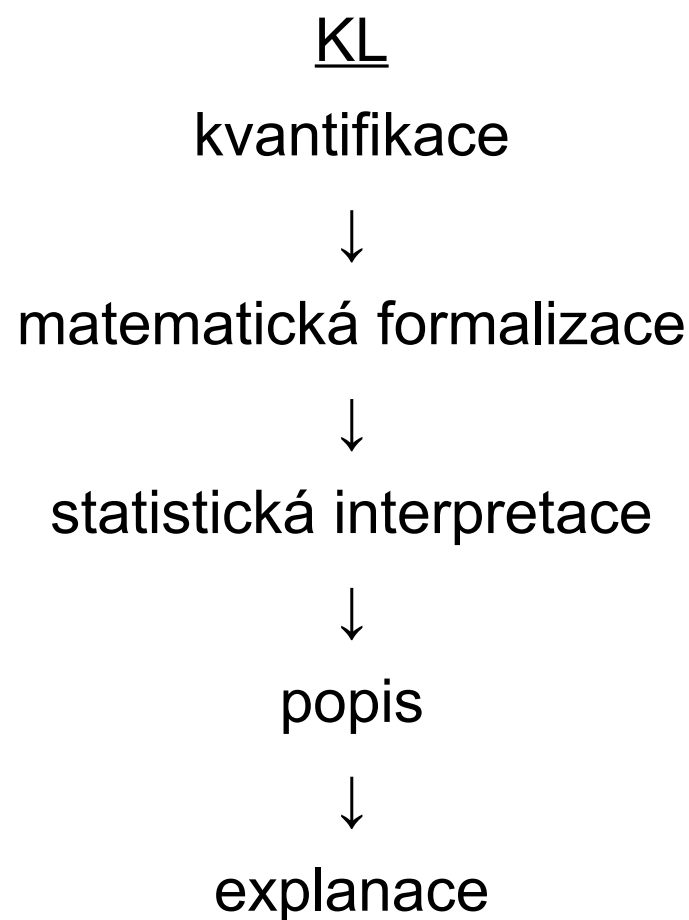




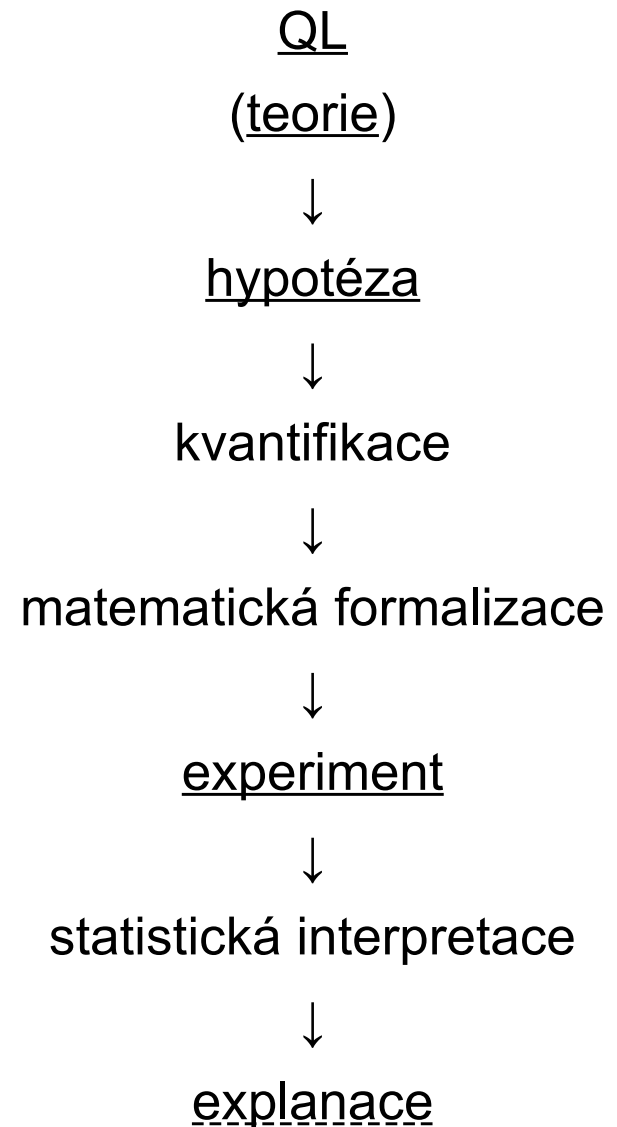
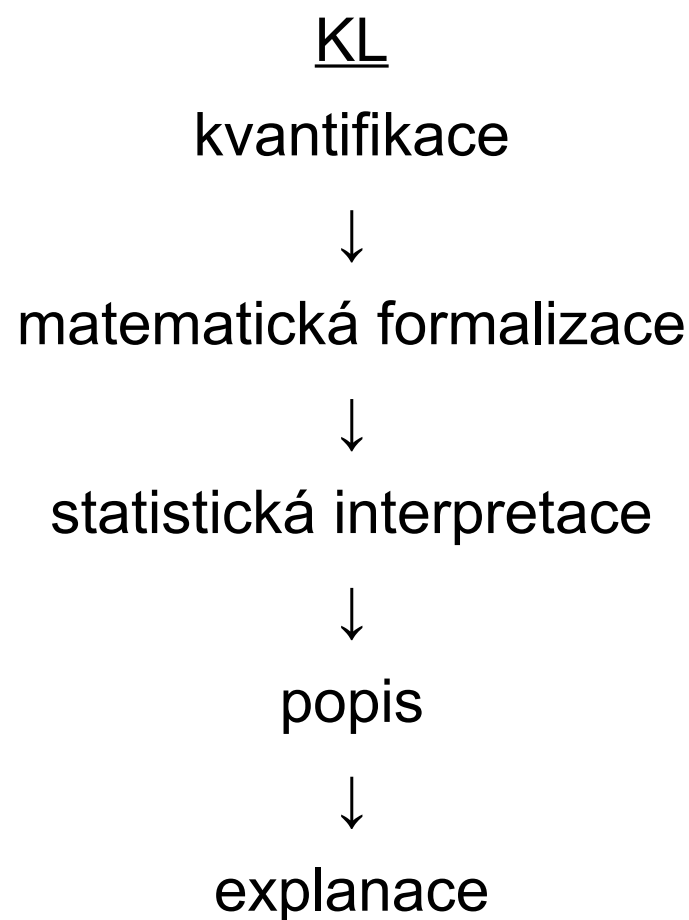
KW vs. TK



Korpusová vs. kvantitativní lingvistika



Korpusová vs. kvantitativní lingvistika



Děkuji za pozornost!

(www.cechradek.cz)