

Začínáme s Bonitem 2 //Word Sketch Engine//

Východisko

Bonito 2 je internetový program, který lze použít na zpracování korpusu libovolného jazyka, je-li tento korpus označován vhodným způsobem.

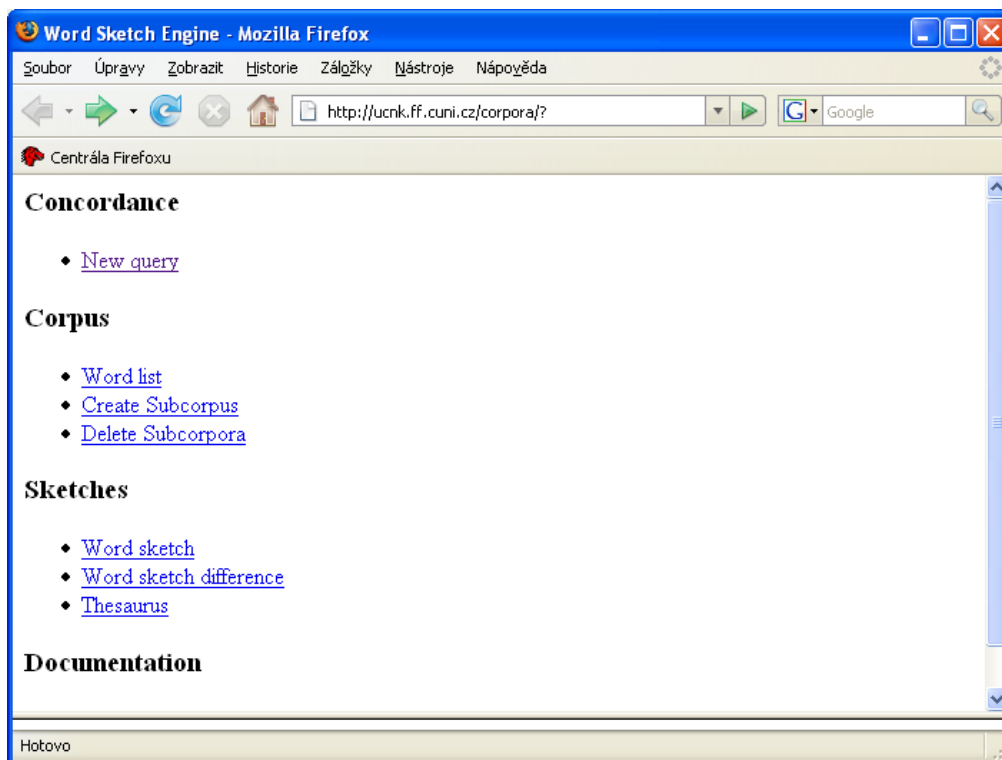
Bonito 2 má řadu funkcí, z nichž základní jsou:

konkordancer (velmi rychlý a vysoce funkční)
program Word Sketch (viz dále).

Více informací o programu Word Sketch naleznete na [Kilgarriff et al 2004 in Proc EURALEX](#).

Domovská stránka

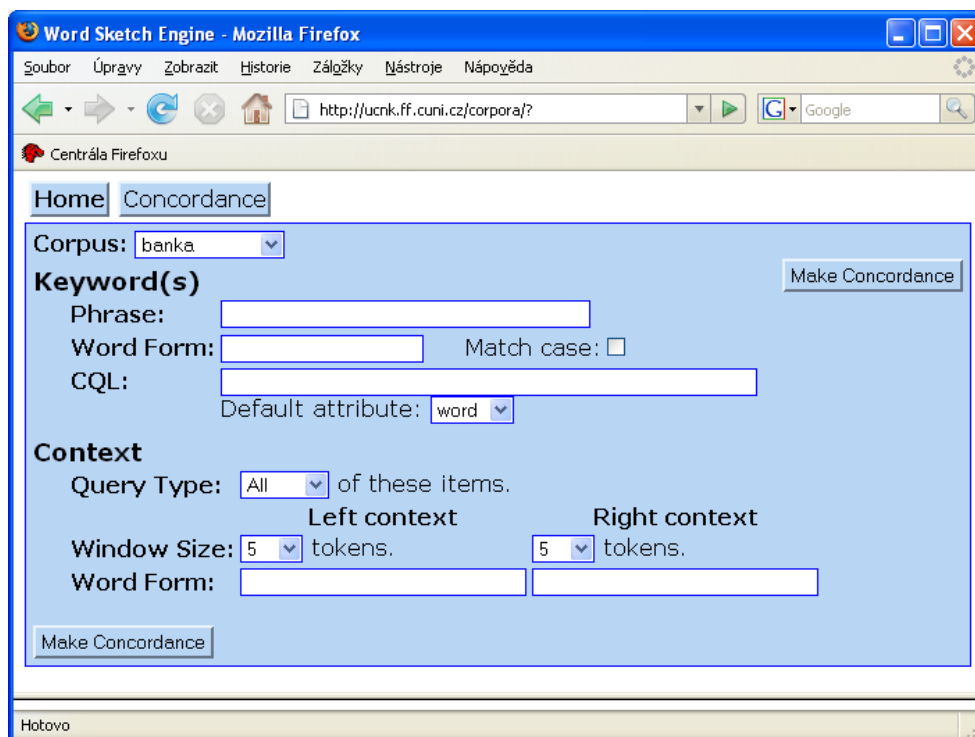
Domovskou stránku jste zobrazili, když jste se do programu přihlásili jako registrovaní uživatelé Českého národního korpusu. Jste tedy na stránce:



Concordance – New query

1.1. Konkordance – Nový dotaz

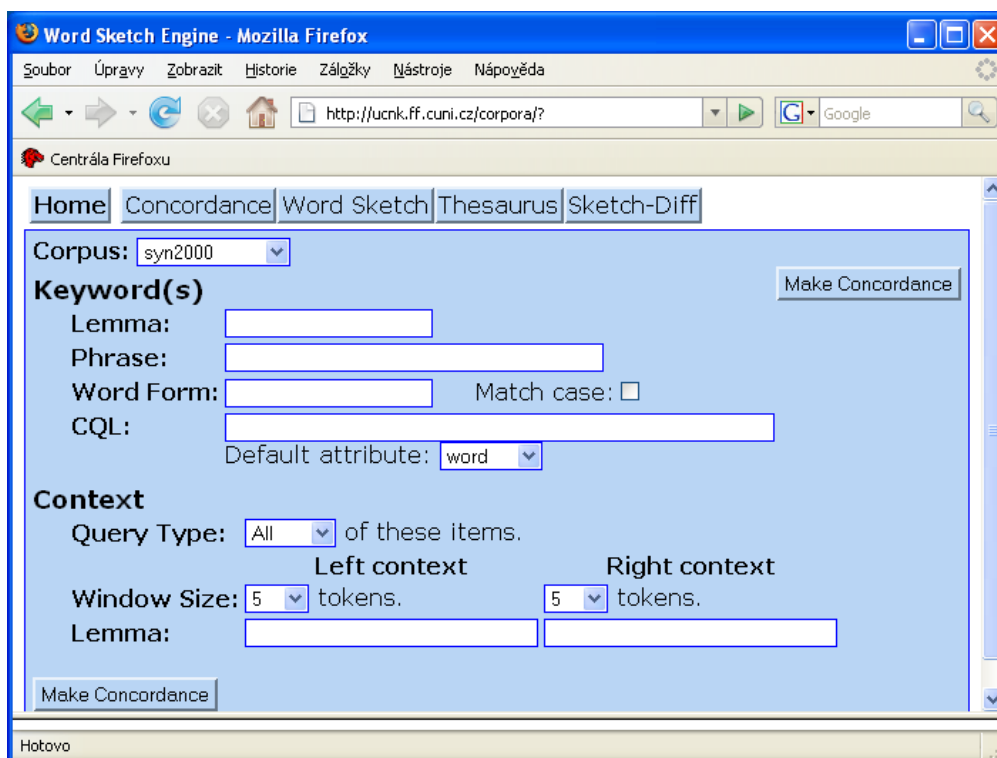
Klikněte na „Concordance - New Query“, dostanete se na stránku zadávacího formuláře.



Nad formulářem je tlačítko pro návrat na domovskou stránku „**Home**“.

Na formuláři je vedle vodícího ukazatele „**Corpus**“ rozevírací roletová nabídka korpusů, se kterými lze aktuálně pracovat. Vytvoříte-li si subkorpus, bude se vám v budoucnu zobrazovat v této nabídce také.

V této nabídce si vybereme korpus SYN2000. Tím se nám stránka změní na:



1.2. Hledané slovo / hledaná slova

Do sekce „**Keyword(s)**“ budeme zadávat námi hledané slovo / hledaná slova.

Za vodící ukazatel „**Lemma**“ zadáváme dotaz, hledáme-li přes lemma.

Za vodící ukazatel „**Word**“ naopak zadáváme dotaz, hledáme-li jednu konkrétní formu, necháme-li přitom prázdné pole „**Match case**“, bude se tvar hledat jak s malými, tak s velkými písmeny (v Bonitu 1 hledání „lc“). Zaklikneme-li možnost „**Match case**“, bude se forma hledat přesně tak, jak jsme ji zadali, tedy např. jen se všemi písmeny malými.

Za vodící ukazatel „**Phrase**“ zadáváme dotaz, hledáme-li kombinace slov, v tomto poli jde však vždy o kombinace sousední. Lze vyhledávat jak na úrovni wordů (je základně nastavena za vodícím ukazatelem „**Default attribute**“, stejně jako je základně nastavena libovolná velikost písmen díky prázdnému poli „**Match case**“), tak na úrovni lemmat. Tuto si můžeme nastavit za vodícím ukazatelem „**Default attribute**“.

Za vodícím ukazatelem „**CQL**“ lze zadávat dotazy pomocí regulárního jazyka. Je to ta forma dotazu, která se nám zobrazila v dotazovacím řádku Bonita 1, vytvořili-li jsme si složitější dotaz pomocí „**Grafického vytváření**“. Nepamatujeme-li si pravidla vytváření dotazu v tomto jazyce, můžeme dotaz vytvořit v „**Grafickém vytváření**“ Bonita 1, dotaz odeslat do korpusu, takže se nám znovu zobrazí v dolním okně. Z tohoto okna pak lze kopírovat.

Za vodícím ukazatelem „**Default attribute**“ si můžeme pro jednoslovné i víceslovné dotazy vybrat, zda chceme hledat na úrovni formy („**Word**“), nebo na úrovni lemmatu („**Lemma**“).

1.3. Kontext

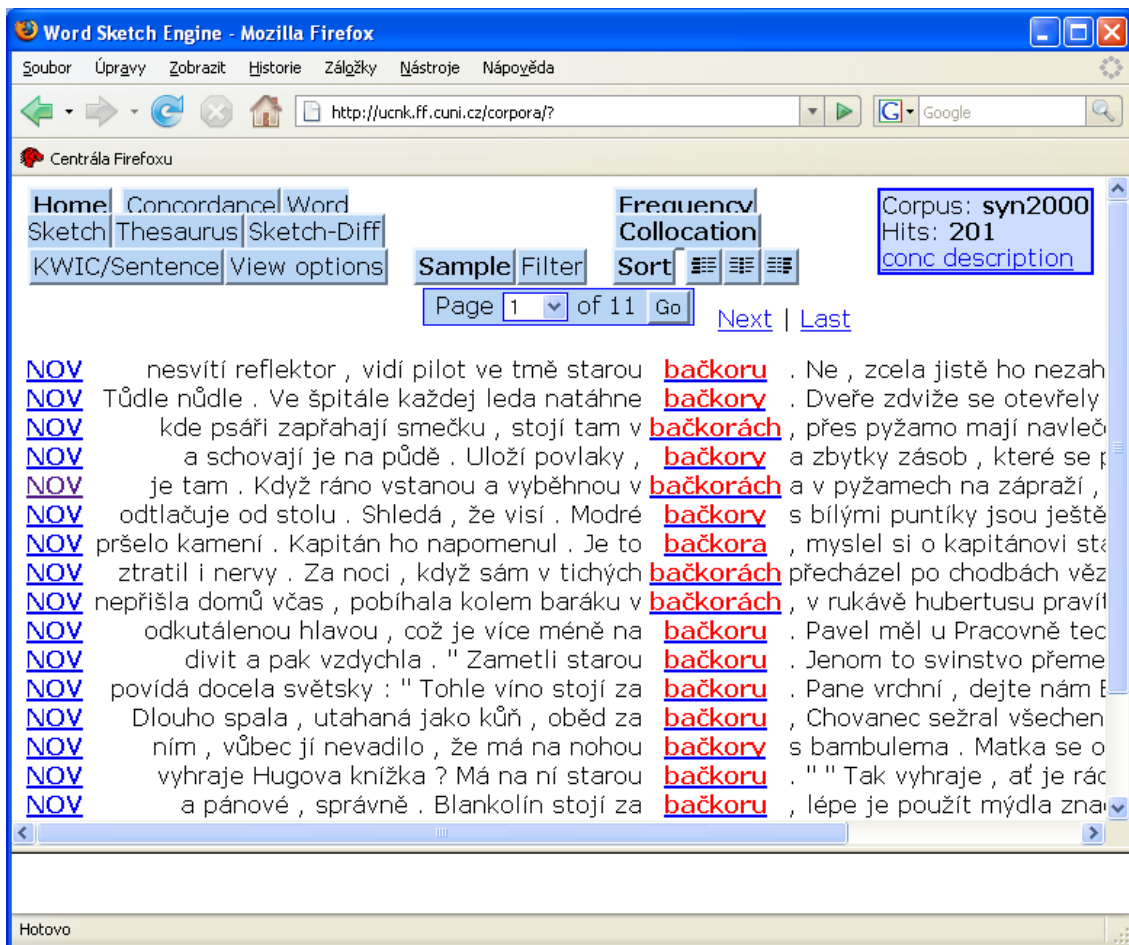
Do sekce „**Context**“ budeme zadávat informace pouze tehdy, jestliže chceme zpřesnit charakteristiku kontextu, ve kterém se námi hledané slovo má nacházet. Jestliže chceme kontext filtrovat.

Za vodícím ukazatelem „**Query Type**“ můžeme vybrat „**All**“, což bude fungovat jako pozitivní filtr (v Bonitu 1 P-filtr), který nám zobrazí pouze ty kontexty klíčového slova, které budou další (NÍŽE) zadané slovo (do vodícího ukazatele „**Lemma**“, jak vidno, zde nelze zvolit úroveň „**Word**“) obsahovat, a to v zadaném rozmezí (defaultně –5 vlevo (Left context) až 5 vpravo (Right context)). Vybereme-li „**None**“, kontexty klíčového slova, které další (NÍŽE) zadané slovo obsahují, se nám nezobrazí (v Bonitu 1 N-filtr).

Za vodícím ukazatelem „**Window Size**“ zadáváme velikost kontextu, ve kterém chceme, aby se další zadané slovo nacházelo. Ale pouze vždy buď vlevo od KWIC, nebo vpravo od něj. Funkce, kterou známe z Bonita 1 – tedy funkce, která mohla pracovat před i za KWIC současně, není v tomto formuláři k dispozici (srovnejte ale 4.4. Filtr !!). Zadáme-li totiž stejné slovo do obou polí (tedy „**Left context**“ i „**Right context**“), hledáme kontext, ve kterém se slovo vyskytlo JAK před, tak současně ZA hledaným klíčovým slovem (a to bývá málokdy). Dotaz na okolí musíme tedy vždy rozdělit na dvě části, levou a pravou (srovnejte ale 4.4. Filtr !!). Na druhou stranu nám tento způsob zpracování umožní relativně jednoduchým způsobem hledat kontext, kde před klíčovým slovem stojí jiné konkrétní slovo, než za ním. V Bonitu 1 jsme to dělali buď postupným použitím dvou P-filtrů za sebou, nebo v Dotaz-Grafické vytváření-Posloupnost, kde jsme mezi hledaná slova vkládali Opakovaně-Libovolná pozice.

1.4. Práce nad vyhledanou konkordancí

Poté, co jsme odeslali formulář do korpusu a co se nám zobrazila hledaná konkordance,



můžeme pokračovat v práci.

V pravém horním rohu vidíme informaci o tom, v jakém se právě nacházíme korpusu („**Corpus**“) a kolik bylo na zadaný dotaz odpovědí, tedy kolik je nalezeno výskytů („**Hits**“). V levé horní části nás tlačítka první řádky odvedou do jiných částí programu, začneme tedy řádkou druhou.

1.4.1. KWIC versus věta

V korpusové lingvistice je vyhledaná konkordance často zobrazována tak, že klíčové slovo (KWIC) je zarovnáno uprostřed, zvýrazněno barvou a kontext ubíhá vlevo a vpravo tzv. do nekonečna. Tlačítko „**KWIC/Sentence**“ nám umožňuje přepínat z tohoto způsobu zobrazení do zobrazení po celých větách (ty jsou pak zarovnány vlevo).

Možnost zobrazit si konkrétní kontext šíře nám ale zůstala, stačí, když levou myší jednou klikneme na KWIC, jehož širší kontext nás zajímá. Objeví se v dolním okně, podobně jako v Bonitu 1.

1.4.2. Možnosti zobrazení

Tlačítko „**View options**“ nás vede k formuláři, na kterém si můžeme vybrat, které hodnoty jednotlivých slov uvnitř textu, tedy „**Attributes**“ chceme zobrazit:

- word = samo slovo / konkrétní forma
- lemma = lemma slova
- tag = celý tag
- lc = slovo / konkrétní forma – ovšem nehledě na velikost mísmen

pos = tag slovního druhu
a některé značky duplicitně ještě jednou či samostatně (tyto jsou zavedeny kvůli zpracování Word Sketch-ů, toto zpracování totiž s korpusovými tagy pracovat nedokáže)

k = značka slovního druhu
g = jmenný rod
c = pád

Můžeme si tu ovšem vybrat pro zobrazení také hodnoty „**References**“, které chápeme jako vnější textu, stojící mimo něj a vztahující se k němu jako k celku:

token number = číselné vyjádření pozice slova v korpusu

doc. type = typ textu
doc. temp = rok vydání
doc. opus = značka jednoznačně identifikující text (a zařazená v seznamu textů)

Relativně málo užívané je zobrazení „**Structures**“, tedy strukturních značek textu (začátky vět, kapitol apod.):

doc = dokument (někdy kapitola románu, někdy celý román, v novinách většinou článek)
s = věta

Tyto možnosti se buď zobrazují jen pro klíčové slovo („**Display attributes**“ „**KWIC tokens only**“), nebo pro všechna slova („**For each token**“).

1.4.3. Vzorek

Tlačítko „**Sample**“ nás vede k formuláři pro „**Random sample**“, tedy náhodný výběr. Volíme zde množství řádek tohoto náhodného vzorku.

1.4.4. Filtr

S filtrem aplikovaným přímo při prvním zadávání dotazu jsme se setkali už v „3.2. Kontext“ a můžeme ho vidět na 3 obrázku v sekci „Context“.

Tlačítko „**Filter**“ nás vede k (v pořadí vlastně druhému) filtru, tento je aplikovaný po vyhledání konkordance; je totožný s funkcemi, které známe z Bonita 1: P-filtr, zde „Filter positive“, a N-filtr, zde „Filter negative“. Pro velikost kontextu, ve kterém chceme filtr použít, platí stejná pravidla jako v Bonitu 1, levý kontext se zapisuje v záporných číslech (tedy např. -3, 0 = filtr bude aplikován do tří pozic vlevo od klíčového slova).

1.4.5. Jednoduché třídění – víceúrovňové třídění

Tlačítko „**Sort**“ nás vede k formuláři pro jednoduché třídění – „**Simple Sort**“ – které známe už z Bonita 1. Nenabízí sice možnost setřídít také klíčové slovo, což Bonito 1 umožňuje, zbylé funkce jsou však shodné. (Setřídění klíčového slova zleva je možné pomocí prostřední třídící ikony ve 4.6. – ovšem možnost třídění klíčového slova zprava /retrográdně/ Bonito 2 nenabízí vůbec).

„**Multilevel sort**“ je funkce, kterou Bonito 1 nemělo. Umožňuje totiž třídít „AŽ“ podle druhého, resp. třetího slova a slova bližší nebrat v úvahu. Podobně umožňuje také volit v rámci těchto pozic POSTUP třídění.

1.4.6. Jednoduché třídění ikonami

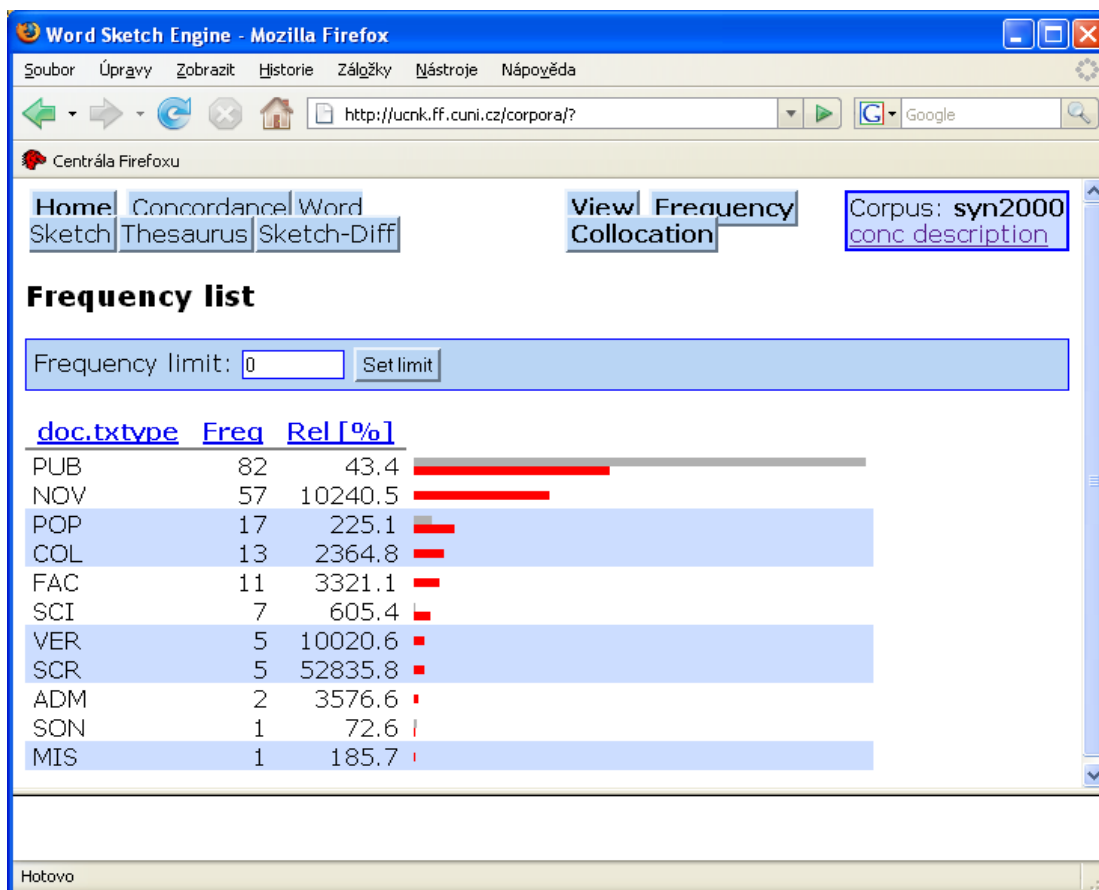
Následující tři tlačítka jsou ikonami pro rychle setřídění – vždy zleva: levého kontextu, klíčového slova, pravého kontextu.

1.5. Frekvence

Tlačítko „**Frequency**“ nás dovede k formuláři, který známe z Bonita 1 jako formulář pro „Frekvenční distribuci“, zde má název „**Multilevel frequency distribution**“. Na rozdíl od Bonita 1 nemůžeme jít dále za pozici 3 vlevo a v pravo od klíčového slova. Nemůžeme také sledovat frekvenci všech atributů, které jsme mohli sledovat v Bonitu 1 – můžeme sledovat: wordu / lc, lemmatu, tagu / pos, a dále: tstype, temp, opus.

Zajímá-li nás víceslovná jednotka, v Bonitu 2 (na rozdíl od Bonita 1) nemůžeme vytvářet statistiky uvnitř tohoto řetězce, protože se zaměřuje jen na okolní kontext. Z klíčového řetězce dokáže statisticky zpracovat pouze jeho první pozici.

Ve formuláři „**Text Type frequency distribution**“ si můžeme vybrat jeden z možných pohledů na distribuci klíčového slova: podle typů textů, podle roku vydání, podle jednotlivých děl. Poměr mezi absolutní frekvencí (zde „Freq“), a relativní frekvencí (tedy frekvencí přepočtenou s ohledem na poměrnou velikost kategorie zde „Rel“) v konkrétním námi zvoleném pohledu je tu znázorněn i sloupcovým grafem.



Sloupcový graf je konstruován následovně:

položky jsou řazeny za sebou podle klesající absolutní frekvence („Freq“). V celé tabulce sloupcového grafu se v souvislosti s tím proporcionalně mění jen zobrazení ČERVENÉ absolutní frekvence.

ŠEDÁ relativní frekvence (přepočtená na základě poměrné velikosti kategorie) je znázorněna proporcionalně pouze vzhledem k sobě odpovídající červené hodnotě absolutní frekvence. Zajímavou novinkou oproti Bonitu 1 je tu možnost nechat si zobrazený seznam „přetřídit“ (abecedně podle kolokujících slov / lemmat, nebo podle relativní frekvence), k tomu slouží aktivní horní popisky sloupců, na které stačí kliknout.

1.6. Kolokace

Tlačítko „**Collocation**“ nás vede k formuláři „Collocation candidates“, který odpovídá v Bonitu 1 funkci Statistika-Kolokace. I zde si vybíráme, na čem chceme statistiku počítat (Attribute: word – lemma – tag – lc – pos – k – g – c), v jak velkém kontextu se bude kolokace zjišťovat („In the range from“). Jaká má být minimální frekvence započítávané jednotky v korpusu („Minimum frequency in corpus“) a jaká má být minimální frekvence započítávané jednotky v námi stanovené velikosti kontextu („Minimum frequency in given range“). Vodící ukazatel „Maximum number of displayed lines“ poukazuje na možnost zvolit si velikost seznamu, který se z korpusu spočítá.

Podobně jako v Bonitu 1 mají zobrazené sloupce scorů aktivní horní popisky sloupců, takže chceme-li získaný seznam přetřídit podle jiného scoru, stačí na daný horní popisek kliknout. Na rozdíl od Bonita 1 jedna je tu ale 6 statistických měř, podle kterých můžeme nechat kolokace setřídit („**Sort by**“), a rozdíl je i v tom, že si můžeme vybrat, které z možných měř se nám ve výsledné tabulce zobrazí („**Show functions**“).

Na rozdíl od Bonita 1 tu máme velmi zajímavou možnost jednoduše se na kolokáty ze zobrazeného seznamu podívat přímo do korpusu. Zcela vlevo máme funkce **p/n**, které jsou shodné s P-filtrem a N-filtrem Bonita 1. Klikneme-li tedy na P, zobrazí se kolokace daného slova (ze seznamu) s naším klíčovým slovem a to podle hodnot, které jsme si zadali do formuláře Kolokace (tedy např. pokud je v rozmezí, které jsme zadali, pokud je v korpusu a v okolí KWICu v počtech, které jsme zadali).

	<u>Freq</u>	<u>T-score</u>	<u>MI</u>	<u>MI3</u>	<u>log likelihood</u>	<u>min. sensitivity</u>	<u>salience</u>
p/n Kožetvorby	3	1.732	18.198	21.368	71.571	0.015	25.228
p/n fajfka	3	1.732	13.795	16.964	51.487	0.015	19.123
p/n bačkora	4	2.000	13.547	17.547	67.283	0.020	21.803
p/n pyžamo	6	2.449	13.336	18.506	99.213	0.017	25.951
p/n kostkovaný	6	2.449	13.272	18.441	98.670	0.016	25.825
p/n župan	6	2.449	13.136	18.306	97.533	0.015	25.561
p/n about	3	1.732	12.521	15.691	46.149	0.010	17.358
p/n plyšový	3	1.732	12.378	15.548	45.552	0.009	17.160
p/n zaklepat	3	1.731	11.452	14.622	41.690	0.005	15.876
p/n ponožka	3	1.731	11.120	14.290	40.304	0.004	15.415
p/n natáhnout	6	2.448	11.092	16.261	80.465	0.004	21.583
p/n teplý	7	2.642	9.515	15.130	78.611	0.001	19.787
p/n bota	4	1.996	8.891	12.891	41.402	0.001	14.310
p/n postel	4	1.994	8.390	12.390	38.628	0.001	13.503
p/n pivo	3	1.723	7.598	10.768	25.677	0.000	10.533

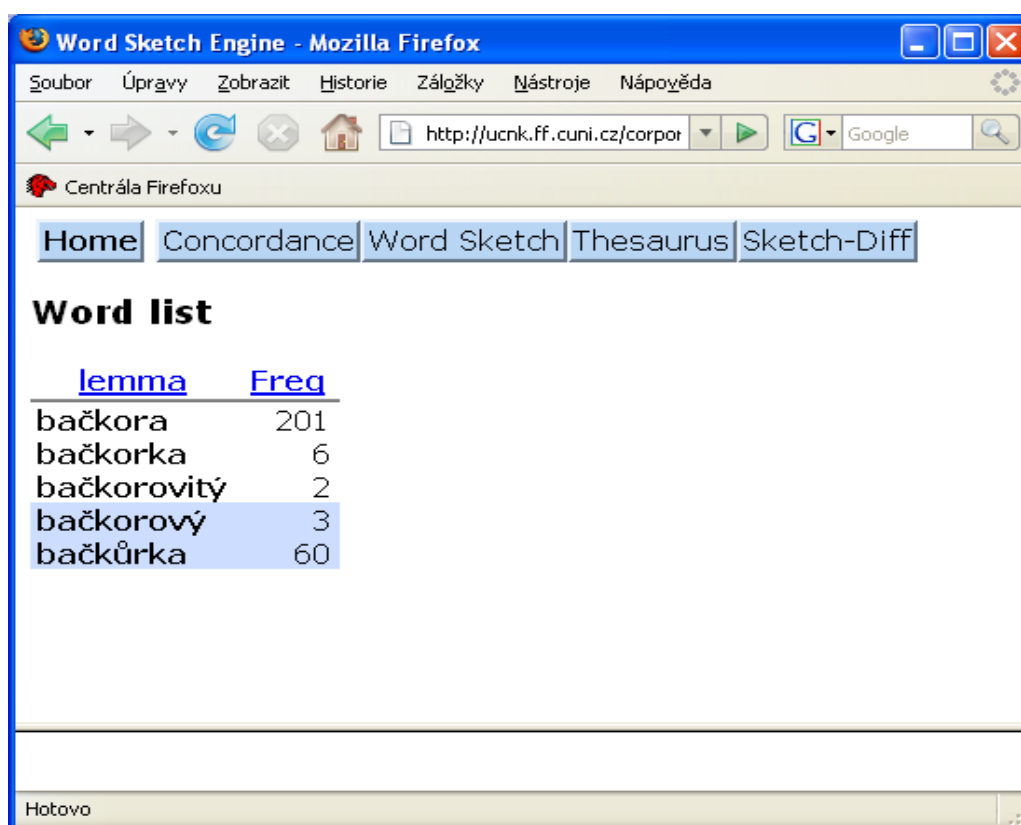
Corpus

2.1. Seznam slov

Pod vodícím ukazatelem "**Word list**" se skrývá funkce, kterou známe z Bonita 1 jako Seznam slov. V Bonitu dvě je to funkce vyčleněná ze sady funkcí "Concordance" ven – musíme tedy znovu zadat korpus, ve kterém chceme pracovat. Podobně jako v Bonitu 1 můžeme zadávat úroveň, kterou prohledáváme ("Attribute": word, lc, lemma, tag, pos, k, g, c), minimální frekvenci hledané jednotky ("Minimum frequency in corpus") a maximální počet zobrazených jednotek ("Maximum number of displayed lines"). Zobrazený výsledek můžeme – podobně jako v Bonitu 1 – přetřídit podle frekvence (kliknutím na "Freq"), nebo zpět podle abecedy (kliknutím na "word").

Na rozdíl od Bonita 1 ale nefunguje druhé kliknutí. V Bonitu 1 totiž toto druhé kliknutí setřídilo kategorii "zpětně" – tedy od nejnižší frekvence k nevyšší, od konce abecedy k začátku (ovšem i tentokrát prohrádně).

Na rozdíl od Bonita 1 také nelze v tomto zobrazení vybrat jednu položku a nechat si zobrazit její korpusové konkordance.



2.2. Vytvoření subkorpusu

Vodící ukazatel "**Create subcorpus**" známe z Bonita 1 jako Vytvoření subkorpusu. V Bonitu 2 klikneme na ukazatel a ve formuláři jako první krok zadáme korpus, ze kterého chceme subkorpus vytvářet. Pak už jen vyplníme jméno subkorpusu ("New subcorpus name") a zaškrtneme v kategoriích, ve kterých chceme vytvářet

- typy textu = doc.txttype
- rok = doc.temp
- konkrétní dílo = doc.opus

konkrétní hodnoty. V pravém sloupci vidíme velikost dané kategorie, vyjádřenou množstvím slov.

2.3. Zrušení subkorpusu

Jednoduchý a názorný formulář se nám objeví pod "**Delete Subcorpora**".

Sketches

3.1. Nárys užívání

I v češtině často používáme nepřeložené označení "**Word sketch**", které zastupuje relativně komplexní funkci, která spojuje následující úkony:

- 1) spočítání absolutní frekvence kolokace
- 2) spočítání hodnoty salience (MI-score krát logaritmus frekvence = tedy nové MI-score, které je upraveno tak, aby se velmi nízké frekvence poněkud znormalizovali) a seřídění podle této hodnoty
- 3) rozdělení seznamu podle syntaktických vztahů (uvnitř kterých zůstává seřídění podle hodnoty salience) /seznam následuje/.

Podobně jako u Seznamu slov, musíme i zde na začátku zadat korpus, ve kterém chceme pracovat. Dále zadáme "Lemma", které nás bude zajímat, možnost pracovat na úrovni "Word" tu nemáme.

"Sort grammatical relations" bylo snad původně zamýšleno jako funkce, která – pokud ji NEzvolíme – spočítá hodnoty kolokací, ale nebude je dělit podle syntaktických vztahů. Funkce není aktivní a pokud chceme mít nedělený seznam seříděný podle salience, použijeme v "Concordance" po vyhledání klíčového slova funkci "Collocation".

"Minimum frequency" je zde myšleno jako minimální frekvence v okolí KWIC-u, se kterou se bude ještě pracovat, nikoli tedy v minimální frekvence v celém korpusu (tu zadávat nemůžeme, na rozdíl od Bonita 1).

"Minimum salience" je minimální hodnota salience, která se bude ve výsledcích zobrazovat.

"Maximum number of items in a grammatical relation" je maximální počet jednotek zobrazených v jedné kategorii (syntaktických vztahů).

Rozmezi (velikost okna), ve kterém jsou kolokace počítány, je -5 / 5 a nelze bohužel změnit.

Nejčastěji zobrazované syntaktické kategorie jsou:

<u>a modifier</u>	přívlastek (shodný) stojící vlevo (od KWIC) panelový dům
<u>prec X</u>	X = konkrétní předložka, např. před
<u>prec "před"</u>	předložkový vztah se substantivem vlevo (od KWIC) stejně jmenný jako slovesný trávník před domem vyběhnout před dům
<u>post X</u>	X = konkrétní předložka, např. s
<u>post "s"</u>	předložkový vztah se substantivem vpravo (od KWIC) stejně jmenný jako slovesný dům se zahradnou

<u>gen 1</u>	genitivní vztah 2 substantiv (KWIC je v nominativu) Dům kultury ...
<u>gen 2</u>	genitivní vztah 2 substantiv (KWIC je v genitivu) majitel domu
<u>prec verb</u>	verbum vlevo stavějí domy
<u>post verb</u>	verbum vpravo dům chátrá
<u>post inf</u>	infinitiv vpravo dům prodat se běží podívat
<u>byt adj</u>	vztah s pomocným slovesem "být" a predikátovým adjektivem (prvkem tagovaným jako adjektivum) dům je prázdný
<u>prec prep</u>	předložka vlevo před domem
<u>post prep</u>	předložka vpravo běžel ke (dvěřím)
<u>has subj</u>	"má podmět (jaký)" vztah verbálního KWIC-u a substantiva v nominativu lhůta běží
<u>is subj of</u>	"je podmětem pro" vztah KWIC-u v nominativu a verba dům vyhořel
<u>is obj2 of</u>	"je předmětem ve 2 pádě pro" vztah KWIC-u v genitivu a verba se netýká domu
<u>is obj3 of</u>	"je předmětem ve 3 pádě pro" vztah KWIC-u v dativu a verba šel domu !! chybně tagováno
<u>has obj4</u>	"má (tento) předmět ve 4 pádě" vztah verbálního KWIC-u a substantiva v akuzativu běžel maraton
<u>is obj4 of</u>	"je předmětem ve 4 pádě pro" vztah KWIC-u v akuzativu a verba prohledali dům
<u>has obj7</u>	"má (tento) předmět v 7 pádě" vztah verbálního KWIC-u a substantiva v instrumentálu běží ulicí
<u>is obj7 of</u>	"je předmětem v 7 pádě pro" vztah KWIC-u v instrumentálu a verba prochází domem
<u>coord</u>	vztah lexikálně vyjádřené koordinace s jednotkou stejné třídy dům a zahradu se otočil a běžel

Výhodnou a jistě často užívanou funkcí bude možnost okamžitě se podívat na konkrétní kolokát, který nás zaujal. stačí, když klikneme na aktivní absolutní frekvenci kolokace

Word Sketch Engine - Mozilla Firefox

Soubor Úpravy Zobrazit Historie Záložky Nástroje nápověda

http://ucnk.ff.cuni.cz/corpora/?

Centrála Firefoxu

Home Concordance Word Sketch Thesaurus Sketch-Diff

dům syn2000 freq = 54895 [change options](#)

a modifier	27266	2.0	prec před	635	9.7	post u	767	9.5
panelový	541	62.55	trávník	16	26.41	zvon	209	69.99
čínžovní	264	60.19	chodník	16	25.72	hybern	52	57.74
bílý	1966	57.7	vyjít	26	24.26	matka	137	47.54
obytný	543	57.14	zaparkovat	8	22.99	prsten	35	39.28
azylový	273	55.9	zastavit	20	22.43	Jonáš	11	29.51
besední	159	54.79	prostranství	7	18.43	rytíř	18	28.22
obecní	842	53.88	parkovat	5	17.23	kaňon	10	27.29
kulturní	1417	53.3	lavička	7	16.91	Glaubic	5	26.35
obchodní	1871	53.24	zahradka	6	16.73	štika	10	25.68
rodinný	852	49.88	vyběhnout	6	16.68	loreta	7	25.63
bankovní	633	48.05	auto	12	15.19	jednorožec	8	23.67
lidový	799	46.89	stát	31	13.78	beránek	6	18.86
tančící	112	44.81	parkoviště	5	13.2	orel	7	17.86
Wortnerův	61	44.41	zahrada	7	11.95	labuť	5	16.75
posádkový	85	42.86	ulice	10	11.31	Novák	7	16.15
bytový	412	40.89	čekat	9	10.29	řeka	10	15.17
rodný	225	39.43	sedět	6	10.02	jezero	7	14.85
měšťanský	109	39.25	být	24	8.25	kolo	6	7.63

Hotovo

3.2. Rozdíly v užívání 2 jednotek

Tato funkce, tedy "Word sketch difference" zpracovává způsobem popsáním v 3.1. dvě jednotky, jejich výsledky porovnává a zobrazuje barevně odlišně (první položka je zelená a různý stupeň přináležitosti k ní je naznačen tónem/sytostí barvy, druhá položka je červená, stupeň přináležitosti je řešen stejně). Tento způsob je zřetelně vidět na tabulce v záhlaví výsledků, kde je vyznačeno i střední, nespécifické pásmo krémové barvy a hraniční hodnoty odlišných tónů/sytostí barev.

chata	21	14	7	0	-7	-14	-21	chalupa
--------------	----	----	---	---	----	-----	-----	----------------

Celkové zobrazení je ovšem celkem názorné. Má dvě varianty. V Bonitu 2 je předem zaškrtnuta za vodícím ukazatelem "Separate blocks" položka "common/exclusive blocks".

Toto funkce nám vytvoří zobrazení:

- jedno pole pro společné jednotky
- další dvě pole pro jednotky specifické pro každé z klíčových slov

Word Sketch Engine - Mozilla Firefox

Soubor Úpravy Zobrazit Historie Záložky Nástroje nápověda

http://ucnk.ff.cuni.cz/corpora/index.html

Centrála Firefoxu

chata/chalupa syn2000 freq = 2640/2716 [change options](#)

Common patterns

chata	21	14	7	0	-7	-14	-21	chalupa
--------------	----	----	---	---	----	-----	-----	----------------

prec_na	263	276	5.7	5.8
jezdit	13	16	20.2	22.3
odjet	6	10	15.7	20.9
jet	9	11	15.8	17.6
dovolená	6	7	15.1	16.5
pobyt	8	5	15.8	11.6
víkend	6	5	14.2	12.5
být	18	18	10.1	9.9

prec_prep	757	937	2.4	2.9
na	351	437	25.6	26.6
u	36	30	16.7	13.9
do	68	63	13.8	11.8
před	24	18	13.5	10.2
k	54	34	13.3	8.3
z	39	69	9.1	12.7
kolem	10	10	11.8	11.1
v	110	153	10.0	11.6
od	16	21	8.1	9.1
za	10	27	3.8	9.1
pro	6	9	1.7	2.6

a_modifier	838	724	1.4	1.1
rekreační	111	71	56.5	49.4
horský	115	31	50.0	30.6
zděný	14	12	29.8	28.2
dřevěný	23	25	25.9	27.6
starý	5	46	5.0	23.9
opuštěný	9	6	18.8	15.2
víkendový	8	7	18.0	17.1
malý	12	9	9.3	7.9

prec_do	81	69	4.4	3.6
vloupání	22	6	45.6	25.1
vloupat	6	5	25.1	23.1

post_v	194	161	2.9	2.3
---------------	-----	-----	-----	-----

gen_2	210	162	1.1	0.8
majitel	40	8	34.4	16.1

is_subi_of	76	90	0.5	0.6
-------------------	----	----	-----	-----

Hotovo

Protože vedle pole pro jednotky společné jsou zde další dvě pole, máme dvojí možnost zadat množství zobrazovaných jednotek: jednak v poli společném "Maximum number...of the common block", jednak v polích disjunktích "Maximum number...of the exclusive block".

Pokud ovšem za vodícím ukazatelem "Separate blocks" zvolíme položku "all in one block", pak budete mít všechny jednotky statistikou považované za relevantní v jediném společném přehledu.

Pozor!!!

Barva tu nesouvisí s výlučností pro jedno z klíčových slov, souvisí pouze s hodnotou salience, která může být i pro slova (v dané dvojici) specifická pro jedno z klíčových slov relativně nízká.

Srovnajte v tabulce níže např. kolokát „oblast“.

chata/chalupa syn2000 freq = 2640/2716 [change options](#)

chata	21	14	7	0	-7	-14	-21	chalupa
--------------	----	----	---	---	----	-----	-----	----------------

prec_u	35	25	9.5	6.6
kurt	6	0	24.4	0.0

prec_na	263	276	5.7	5.8
jezdit	13	16	20.2	22.3
odjet	6	10	15.7	20.9
jet	9	11	15.8	17.6
dovolená	6	7	15.1	16.5
pobyt	8	5	15.8	11.6
wikend	6	5	14.2	12.5
být	18	18	10.1	9.9
rodina	5	0	9.4	0.0

prec_do	81	69	4.4	3.6
vloupání	22	6	45.6	25.1
vloupat	6	5	25.1	23.1

prec_k	48	26	3.4	1.8
dojít	5	0	14.0	0.0

post_v	194	161	2.9	2.3
Seč	7	0	29.6	0.0
Krkonoše	7	7	23.3	23.9
Čechy	0	11	0.0	19.8
osada	6	0	18.4	0.0
obec	6	11	11.9	18.1
okolí	7	0	15.6	0.0
hora	5	7	11.8	15.3
les	7	0	14.9	0.0
oblast	8	0	12.1	0.0

prec_prep	757	937	2.4	2.9
na	351	437	25.6	26.6
u	36	30	16.7	13.9
do	68	63	13.8	11.8
před	24	18	13.5	10.2
k	54	34	13.3	8.3
z	39	69	9.1	12.7
kolem	10	10	11.8	11.1
v	110	153	10.0	11.6
nad	0	14	0.0	10.7
od	16	21	8.1	9.1
za	10	27	3.8	9.1
pod	7	0	6.6	0.0
po	0	12	0.0	4.5
pro	6	9	1.7	2.6
o	7	0	0.4	0.0
s	0	10	0.0	0.2

prec_z	33	61	1.5	2.8
vyhnat	0	7	0.0	25.0
vyjít	0	5	0.0	15.0

3.3. Tezaurus

Za vodícím ukazatelem „**Thesaurus**“ najdeme formulář pro hledání lemmat (v konkrétním korpusu, který si musíme nejdříve vybrat), v jejichž okolí se vyskytuje určitý počet stejných jednotek jako v okolí zadaného klíčového slova. Jsou tedy porovnány výsledky „Sketchů“ pro celý korpus a výsledný seznam obsahuje ty jednotky, u kterých byla shledána podobnost s klíčovým slovem („Minimum similarity between cluster items“) vyšší, než ve formuláři zadaná (nebo jí rovná).

Jak je tato „podobnost“ počítána, není jasné. Z výsledků není zřejmé bohužel ani to, k čemu referují hodnoty zobrazované za každým uváděným lemmatem. Lze jen říci, že se výsledky skutečně proměňují v závislosti na

- počtu zobrazovaných lemmat („Maximum number of items“)
- minimální podobnosti s klíčovým slovem („Minimum similarity between cluster items“)

Za zmínku zde stojí organizace (shlukování) zobrazovaných lemmat. Lemmata nejsou zobrazována podle jim odpovídajících hodnot (jejichž referenci bohužel neznáme), ale podle vzájemných podobností. V tezauru pro slovo „dům“ tedy můžeme vidět následující shluky – vybíráme:

budova 0.382	byt 0.361	objekt 0.295	prostor 0.201	stavba 0.2	zařízení 0.182	areál 0.162
město 0.261	obec 0.194	země 0.188	Praha 0.157	republika 0.155		
místnost 0.191	pokoj 0.172	hala 0.166	sál 0.152			
obchod 0.175	služba 0.165					
rodina 0.175	dítě 0.154					
majetek 0.161	pozemek 0.151					