

PLIN033_3

Přegenerování a podgenerování –
dva problémy automatické analýzy
přirozeného jazyka, konkrétně
slovotvorby

Přegenerování

- *Formální definici (algoritmu) odpovídají jednotky, které tvoří homogenní skupinu (tu, kterou se prostřednictvím formálního zadání snažíme definovat), ale i jednotky, které jsou vůči této skupině heterogenní. Tento jev spadá na vrub obecné vlastnosti přirozeného jazyka, jíž je víceznačnost (homonymie) na všech úrovních.*

Podgenerování

- *Rubem téže mince je tzv. podgenerování, tedy případ, kdy formální zadání je vymezeno příliš úzce, takže nejsou zachyceny jednotky, které se jeho prostřednictvím snažíme definovat.*

Příklady přegenerování z minulých cvičení

- *Náboženství, nádeničení, ...*
- *Klíč, míč, ...*

Pomocí nástroje *Deriv* a *Morfio* vyhledejte kandidáty na činitelská jména na *-tel*

- Maskulina životná s koncovým řetězcem *tel*

Hledání

značka RE

Deriv

- seznam

1	<input type="checkbox"/>	<u>badatel</u>
2	<input type="checkbox"/>	<u>bavitel</u>
3	<input type="checkbox"/>	<u>bičovatel</u>
4	<input type="checkbox"/>	<u>bořitel</u>
5	<input type="checkbox"/>	<u>bouratel</u>
6	<input type="checkbox"/>	<u>branitel</u>
7	<input type="checkbox"/>	<u>brojitel</u>
8	<input type="checkbox"/>	<u>brzditel</u>
9	<input type="checkbox"/>	<u>buditel</u>
10	<input type="checkbox"/>	<u>budovatel</u>
11	<input type="checkbox"/>	<u>burcovatel</u>
12	<input type="checkbox"/>	<u>bydlitel</u>
13	<input type="checkbox"/>	<u>cenitel</u>
14	<input type="checkbox"/>	<u>cestovatel</u>

Deriv hledání dvojic

- **t\$/k5.*mF>tel/k1gMnSc1**

1: značka RE

2: značka RE odvodit od 1. slova: nahrazované nahrazující

Deriv

- seznam

1	<input type="checkbox"/>	<u>badat, badatel</u>
2	<input type="checkbox"/>	<u>bavit, bavitel</u>
3	<input type="checkbox"/>	<u>bičovat, bičovatel</u>
4	<input type="checkbox"/>	<u>bořit, bořitel</u>
5	<input type="checkbox"/>	<u>bourat, bouratel</u>
6	<input type="checkbox"/>	<u>brojit, brojitel</u>
7	<input type="checkbox"/>	<u>brzdit, brzditel</u>
8	<input type="checkbox"/>	<u>budit, buditel</u>
9	<input type="checkbox"/>	<u>budovat, budovatel</u>
10	<input type="checkbox"/>	<u>burcovat, burcovatel</u>
11	<input type="checkbox"/>	<u>bydlit, bydlitel</u>
12	<input type="checkbox"/>	<u>cenit, cenitel</u>
13	<input type="checkbox"/>	<u>cestovat, cestovatel</u>
14	<input type="checkbox"/>	<u>cvičit, cvičitel</u>
15	<input type="checkbox"/>	<u>čekat, čekatel</u>
16	<input type="checkbox"/>	<u>činit, činitel</u>
17	<input type="checkbox"/>	<u>čistit, čistitel</u>

Morfio

- Seznam

	vzor 1 ▲▼ (fq ▲▼)	vzor 2 ▲▼ (fq ▲▼)
1	žít (3800)	živitel (142)
2	zvěstovat (250)	zvěstovatel (49)
3	zůstavit (38)	zůstavitel (295)
4	zřizovat (555)	zřizovatel (599)
5	zprostředkovat (910)	zprostředkovatel (605)
6	zpracovat (2947)	zpracovatel (460)
7	zplodit (485)	zploditel (11)
8	znečišťovat (143)	znečišťovatel (147)
9	zmocnit (2601)	zmocnitel (29)
10	zhotovit (889)	zhotovitel (776)
11	zásobit (204)	zásobitel (15)
12	zasílat (719)	zasílatel (51)
13	zapisovat (1400)	zapisovatel (85)
14	zajistit (12412)	zajistitel (17)
15	vyzývat (1610)	vyzývatel (51)
16	vyšetřovat (2548)	vyšetřovatel (2394)
17	vystavovat (2286)	vystavovatel (477)
18	vysílat (3721)	vysílatel (40)
19	vypisovat (351)	vypisovatel (16)
20	vykořisťovat (86)	vykořisťovatel (59)
21	vychovat (1091)	vychovatel (389)
22	vyhlašovat (834)	vyhlašovatel (107)
23	vydražit (227)	vydražitel (34)
24	volit (5757)	volitel (117)
25	vnímat (10105)	vnímatel (43)
26	věznit (520)	věznitel (148)
27	věřit (29642)	věřitel (2047)

Přegenerování

- Přít/přítel

341	<input type="checkbox"/>	přisvědčovat , přisvědčovatel
342	<input type="checkbox"/>	přít , nepřítel
343	<input type="checkbox"/>	přít , přítel
344	<input type="checkbox"/>	přítakat , přítakatel

57	rušit (3234)	rušitel (41)
58	ručit (658)	ručitel (197)
59	rozhodovat (7902)	rozhodovatel (79)
60	přít (998)	přítel (23968)
61	přihlašovat (131)	přihlašovatel (30)
62	představit (22703)	představitel (6912)

Důvody přegenerování

- Příliš široké formální vymezení
- Nemožnost užšího formálního vymezení

Podgenerování

- Kde jsou slova jako *ředitel*, *uchvatitel*, *šířitel*, *majitel*, *pisatel*, ... ?
- Zahrnutí alternací do vyhledávání jakožto prostředek zúžení definice hledaných jednotek.

Derivační pravidla a výsledky pro derivaci sloveso – dějové jméno na *-tel*

-tel

substituční pravidlo	příklad		dvojice	<u>přegenerování</u>
<u>(?<lí)t\$/k5.*mF</u> >tel/k1gMnSc1 ¹	<i>učit/učitel</i>	716	716	0
et\$/k5.*mF> <u>itel/k1gMnSc1</u>	<i>velet/velitel</i>	20	16	4 ²
<u>([ěí])t\$/k5.*mF>itel/k1gMnSc1</u>	<i>trpět/trpítel</i>	27	22	5 ³
<u>át\$/k5.*mF>atel/k1gMnSc1</u>	<i>znát/znatel</i>	2	1	1 ⁴
<u>át\$/k5.*mF>ítel/k1gMnSc1</u>	<i>přát/přítel</i>	2 ⁵	2	0
<u>á(.)[iě]t\$/k5.*mF>a\$1itel/k1gMnSc1</u>	<i>uchvátit/uchvatitel</i>	19	18	1 ⁶

Derivační pravidla a výsledky pro derivaci sloveso – dějové jméno na *-tel*

$i(.)it/k5. *mF > i\$1itel/k1gMnSc1$	<i>šít/šítel</i>	7	7	0
$i(.)it/k5. *mF > e\$1itel/k1gMnSc1$	<i>řít/řítel</i>	2	2	0
$ou(.)[ie]t/k5. *mF > u\$1itel/k1gMnSc1$	<i>vykoupit/vykupitel</i>	11	9	2 ¹⁶⁷
$á(. ch)at/k5. *mF > a\$1atel/k1gMnS1$	<i>dodávat/dodavatel</i>	57	57	0
$i(.)at/k5. *mF > i\$1atel/k1gMnSc1$	<i>užívat/uživatel</i>	15	15	0
$i(.)at/k5. *mF > ě\$1atel/k1gMnSc1$	<i>přispívat/přispěvatel</i>	4	4	0
$y(.)at/k5. *mF > y\$1atel/k1gMnSc1$	<i>obývat/obyvatel</i>	7	7	0
$é(.)at/k5. *mF > e\$1atel/k1gMnSc1$	<i>obléhat/oblehatel</i>	2	1	1 ¹⁶⁸
$ou(.)at/k5. *mF > u\$1atel/k1gMnSc1$	<i>zkoumat/zkumatel</i>	1	1	0

Derivační pravidla a výsledky pro derivaci sloveso – dějové jméno na *-tel*

jmout\$/k5. *mF>jatel/k1gMnSc1 ¹⁶⁹	<i>pronajmout/pronajatel</i>	1	1	0
ýt\$/k5. *mF>ytel/k1gMnSc1	<i>vydobýt/vydobytel</i>	1	1	0
ít\$/k5. *mF>ajitel/k1gMnSc1	<i>mít/majitel</i>	1	1	0
ést\$/k5. *mF>editel/k1gMnSc1	<i>provést/proveditel</i>	1	1	0
ct\$/k5. *mF>žitel/k1gMnSc1	<i>přemoct/přemožitel</i>	1	1	0
^psát\$/k5. *mF> pisatel/k1gMnSc1	<i>psát/pisatel</i>	1	1	0
CELKEM	22 pravidel	898	884	14
Výjimky	<i>neumětel, postihatel, dosažitel, ...</i>			

Přegenerované doklady

- *mučet/mučitel, proset/prositel, těžet/těžitel, zcizet/zcizitel.*
- *mocnět/mocnitel, pět/pitel, zmocnět/zmocnitel, pro-sít/prositel, učit/učitel.*
- *dát/datel.* Tento případ přegenerování by bylo možno eliminovat aplikací podmínky, že substantivum musí být skloňováno podle vzoru *muž*, již by bylo možno zadat v případě, že bychom pracovali se strojovým slovníkem značkováním tak, že by součástí značky byla i informace o flektivním typu (vzor).
- Jedná se o substantiva *přítel* a *nepřítel*. V praxi automatické morfologické analýzy (lemmatizace) nepanuje jednota v interpretaci derivátů se záporkou *ne-* (srv. Osolsobě 2007¹). Řešení tohoto problému přesahuje záměr této práce.
- *zařádit/zařaditel.*
- *boudit/buditel, moučit/mučitel.*
- *ohlédat/ohledatel.*

Přehled alternací

alternace KmV	příklad
<i>o-i</i>	<i>přemoc0t-přemožitel</i>
<i>á-a</i>	<i>psát-pisatel</i>
<i>e-i/ě-i</i>	<i>velet-velitel/ trpět-trpitel</i>
<i>í-i</i>	<i>křtit-křtitel</i>
<i>mou-a</i>	<i>pronajmout-pronajatel</i>
alternace KoV	
<i>á-a</i>	<i>uchvátit-uchvatitel</i>
<i>é-e</i>	<i>obléhat-oblehatel</i>
<i>í-i</i>	<i>šířit-šířitel</i>
<i>í-e</i>	<i>řídit-ředitel</i>
<i>í-ě</i>	<i>přispívat-přispěvatel</i>
<i>ou-u</i>	<i>vykoupit-vykupitel</i>
<i>ý-y</i>	<i>hýbat-hybatel</i>
<i>í-aj</i>	<i>mít-majitele</i>
<i>o-i</i>	<i>p0sát-pisatel</i>
<i>á-í</i>	<i>přát-přítel</i>

Přehled alternací

alternace kf	
<i>c-ž</i>	<i>přemoc0t-přemožitel</i>
<i>h-ž</i>	<i>dosáhnout-dosažitel</i> ¹⁵⁸
konekt	
<i>nou-a</i>	<i>postihnout-postihatel</i>
<i>nou-i</i>	<i>dosáhnout-dosažitel</i>
alternace na švu prefix/kořen (P0e, ki) ¹⁵⁹	

Vyhledávání dvojic

- **at\$/k5.*mF>áč/k1gMnSc1**

1: značka RE

2: značka RE odvodit od 1. slova: nahrazované nahrazující

Přegenerování

Práce s obsahem souboru PLIN033/at_áč

Soubor **PLIN033/at_áč** byl úspěšně uložen (2 s).

Poznámkou je jednotlivý znak, jeden řádek může mít více poznámek.

Jako poznámku nelze použít dvojtečku, čárku a mezeru (budou-li zadány, budou ignorovány).

- | | | |
|----|--------------------------|---|
| 1 | <input type="checkbox"/> | belhat , belháč |
| 2 | <input type="checkbox"/> | brnkat , brnkáč |
| 3 | <input type="checkbox"/> | brumlat , brumláč |
| 4 | <input type="checkbox"/> | bukat , bukáč |
| 5 | <input type="checkbox"/> | česat , česáč |
| 6 | <input type="checkbox"/> | hafat , hafáč |
| 7 | <input type="checkbox"/> | hmatat , hmatáč |
| 8 | <input type="checkbox"/> | chmatat , chmatáč |
| 9 | <input type="checkbox"/> | klepetat , klepetáč |
| 10 | <input type="checkbox"/> | kokrhat , kokrháč |
| 11 | <input type="checkbox"/> | kopat , kopáč |
| 12 | <input type="checkbox"/> | kovat , kováč |
| 13 | <input type="checkbox"/> | krkat , krkáč |
| 14 | <input type="checkbox"/> | orat , oráč |

klepetáč

- Slovník

SSJČ Slovník spisovného jazyka českého

klepetáč

-e m. expr. *velký rak*: chytil k-e za hlavu, aby nestříhl (Pujm.); --> expr. zdrob. **klepetáček**, -čka m. (mn. 1. -čci, -čkové
6. -čcích)

krkáč

- slovník

ssjc Slovník spisovného jazyka českého

***krkáč**

-e m. expr. *lakomec, chamtivec, krkoun* (Herrm., Šlej.)

Důvody

- Polyfunkčnost prostředku (-á-č x –áč)
- Závisí na mimojazykových znalostech
- Obtížně se formálně definuje

Podgenerování

- Nedostatky ve formální definici
- Nepravidelnosti (*vozač, trubač*)
- Jednotky nejsou zachyceny ve slovníku
- Jednotkám nezachyceným ve slovníku chybí interpretace na úrovni lemmatu a morfologické značky

Morfio

- *kout/kouč, klít/klíč, sálat/salač*
- Propria: *máchat/Machač, tykat/Tykač, dědit/Dědič, pískat/Pískač, kopat/Kopač, klapat/Klapač, kovat/Kovač, pleskat/Pleskač, bílit/Bilič*

Typy přegenerování

- hláskové alternace **kořenového vokálu** u derivátů od sloves **III. třídy** podle kmene přítomného (vzor *krýt*)
- hláskové alternace **kořenového vokálu** u ostatních tříd a vzorů
- hláskové alternace **kmenotvorného vokálu** u ostatních tříd a vzorů

Alternace KoV u derivátů sloves podle *krýt*

- *hrát/hráč*
- *chcát/chcáč*
- ? *pít/píč*
- ? *pět/pěč*
- ? *sít/síč*

V korpusech lze najít (SYN)

- *pít* (čaj)/ *čajpíč*
- *žit/žič*
- **!** *šít/šič*

" Po kafi mám blijavku , soudruzi ! Máte čaj . . . ? Já jsem . . . **čajpíč** ! Kdo je tady eště . . . čaj . . . ? "

Reflex, č. 45/2002	, že spadáme do kategorie	žičů /žičů/X@-----	. Jako je někdo optimista
Reflex, č. 45/2002	flegmatik , tak my jsme	žiči /žiči/X@-----	. A když už se
Reflex, č. 45/2002	už se jednou člověk takovým	žičem /žičem/X@-----	narodí , tak jím nikdy

Blesk, 17. 1. 2000	měsíc . " Rekvafikaci na	šiče /šiče/X@-----	mi nabídli , když jsem
Blesk, 17. 1. 2000	Loni bylo na rekvafikační kurs	šiče /šiče/X@-----	, který zajiřtuje zmíněná firma
Blesk, 18. 5. 2000	zájem o práci šiček a	šičů /šičů/X@-----	, kterých tady pracuje sto
Události DNES 20. 9. 2000	uřídí se šiče šiček	šiče /šiče/X@-----	šiček . Těto šiček

A kromě toho u neživotných máme

- *bít/bič*
- *rýt/rýč*

Všimněme si dvojic

- vyprá*á*vět/vypr*a*věč | vyprá*á*věč
- vyjedná*á*vat/vyjedn*a*vač | vyjedná*á*vač

IJP

- <http://prirucka.ujc.cas.cz/?id=730#nadpis14>
- **2 Střídání krátkých a dlouhých samohlásek při tvoření slov**
- Příklady nikoli pravidla (?seznamy výjimek)

Úkol na 29.10. 2014

- Pomocí nástrojů *Deriv* a *morfio* vyhledejte kandidáty na dvojice sloveso-jméno prostředku na -dlo.
- Popište případy přegenerování popř. podgenerování