

# Disambiguace několika českých slov pomocí Brillova taggeru

Přemysl Šrubař

## Specifikace úlohy

Cílem je disambiguace (zjednoznačnění) několika vybraných českých slov pomocí *Brillova taggeru*. Jedná se o tato slova a jejich značky:

se	k3, k7
je	k5, k3
vedoucí	k1gM, k1gF, k2gM, k2gF
Jana	k1gM, k1gF

## Příprava dat

Data byla získána z korpusu *Desam*. Pro každé vybrané slovo (se, je, vedoucí, Jana) v kontextu pěti řetězců (pět před i po vybraném slově), pouze značky (bez lematu).

Pro použití v *Brillově taggeru* bylo třeba získaný korpus upravit:

- Doplnit konkrétní značku za prázdné značky. Například tečka nebo čárka za větou je v korpusu *Desam* označena pouze lomítkem bez značky. Takovéto řetězce byly označeny značkou *ke*.
- Odstranit odsazení a prázdné řádky.
- Odstranit ostré závorky, kterými jsou označena vybraná slova. Např: `<Jana/k1gMnSc2>`.
- Vynechat řádky, které obsahují značky jako `<x>`, `<z>` apod.
- Doplnit značku ke slovům, která nejsou označována. To bývají převážně slovní spojení jako například *Tyranosaurus Rex/k1gMnSc1*. Těmto slovům byla dána značka *ke*.
- Výskyt více lomítek za sebou přepsat na znak |.

Pro zjednodušení byl navíc u značek ponechán pouze druh slova a rod (pro rozlišení slova *Jana*).

Všechny tyto úpravy obstarává skript [desam2brill.prl](#)

## Brillův tagger

Přeložený *tagger* (win32): [Tagger Win32.zip](#).

Naučená data: [SeJeJanaVedouci.zip](#).

Značkování pomocí Brillova *taggeru* probíhá zhruba takto:

- Všem slovům se přiřadí jejich nejčastější značka.
- Pomocí pravidel pro značkování neznámých slov se přiřadí značka i neznámým slovům
- Aplikují se kontextová pravidla – opravují se chyby.

Naučená data jsou v těchto souborech:

<b>Soubor</b>	<b>Popis</b>	<b>Příklad</b>
LR	Lexical-rules - Pravidla pro značkování neznámých slov.  Naučený <i>tagger</i> obsahuje těchto pravidel 1207 (pro angličtinu zhruba 140).	<b>o char k5 5328.27817522851</b> Pokud se ve slově vyskytuje písmenko <i>o</i> , označuj slovo <i>k5</i>
		<b>k1 a fchar k1gF 1257.97859349202</b> Změň značku <i>k1</i> na <i>k1gF</i> , pokud se ve slově vyskytuje <i>a</i> .
		<b>ni hassuf 2 k1gN 987.397887274155</b> Má-li slovo suffix <i>ni</i> , označuj slovo jako <i>k1gN</i> .
		<b>k5 v fgoodright k1gF 31.3566037735849</b> Změň značku <i>k5</i> na <i>k1gF</i> , pokud se napravo vykytuje slovo <i>v</i> .
CLR	Contextual-rules – Kontextová pravidla.  Naučený <i>tagger</i> obsahuje těchto pravidel 448 (pro angličtinu zhruba 250).	<b>k2gF k2gI NEXT1OR2TAG k1gI</b> Změň <i>k2gF</i> na <i>k2gI</i> , pokud má jedno ze dvou následujících slov <i>k1gI</i> .
		<b>k2gI k2gN NEXTTAG k1gN</b> Změň <i>k2gI</i> na <i>k2gN</i> , pokud má následující slovo značku <i>k1gN</i> .
		<b>k3gI k3gM PREV1OR2OR3TAG k1gM</b> Změň <i>k3gI</i> na <i>k3gM</i> , pokud má jedno ze tří předchozích slov <i>k1gI</i> .
		<b>k3 k7 PREV1OR2WD se</b> Změň <i>k3</i> na <i>k7</i> , pokud je jedno ze dvou předchozích slov <i>se</i> .
BGL	Viz dále.	
TL	Viz dále.	

Ve všech souborech je použité kódování češtiny *Win CP 1250*.

Ve vstupním textu musí být všechny řetězce odděleny mezerou, včetně interpunkce. Například:

*Po zápase řekl: "Cítím se dobře."*

musí být:

*Po zápase řekl : " Cítím se dobře . "*

Také by měla být každá věta na samostatném řádku, ale není to nutné.

Učení *taggeru* trvalo přibližně 66 hodin. Z toho 56 hodin trvalo učení modulu pro značkování neznámých slov, které je napsáno v *perlu*. Učení kontextových pravidel je napsáno v jazyku *C*.

Naučený *tagger* si můžete vyzkoušet přes [webové rozhraní](#).

## Proces učení

Soubory nutné k učení (pod *Win32*): [Learn\\_Win32.zip](#).

Většina utilit, které je třeba k učení, jsou napsána v jazyku *Perl*. Je také nutno použít některé základní programy *Unixu* : *cat*, *dos2unix*, *(g)awk*, *ls*, *mv*, *perl*, *rm*, *sort*, *tee*, *(unix2dos)*.

I některé utility napsané v jazyku *C* vyvolávají systémové příkazy. Dokonce i *tagger.exe* vyvolává programy *cat* a *tee*, nastěsí ne při samotném značkování (bez použití volitelných přepínačů).

Pro učení pod *Win32* je proto vhodné nainstalovat *Cygwin* - některé programy *unixu* portované pro *Win32*. Pozor na program *sort*, který je součástí novějších *windows*, ale není kompatibilní s unixovským *sort* (přepínače *-rn*). Bývá většinou v cestě.

Pokud nemáte nebo nechce instalovat *Cygwin*, stačí rozbalit archiv [CygwinCropped.zip](#) do stejného adresáře jako archiv [Learn\\_Win32.zip](#) a přidat soubor *RegCygwin.reg* do registrů (obsahuje pouze tři záznamy o mapování adresářů */usr/bin*, */usr/lib* a */*, které se namapují na aktuální adresář).

Učení *Brillova taggeru* se skládá ze dvou hlavních částí: hledání pravidel pro značkování

neznámých slov a hledání pravidel podle kontextu. K zahájení učení (všech částí) je možné použít skript

```
Learn.bat corpus1 [ [ [ [ corpus2 ] corpus3 ] .. ] corpus9 ]
```

Kde *corpus1 .. corpus9* jsou korpusové soubory ve formátu *desam*. Je-li zadáno více souborů, jsou nejdříve spojeny dohromady.

Dále podrobněji:

### **Učení modulu pro značkování neznámých slov:**

Před zahájením samotného učení je třeba nejdříve udělat několik úprav, případně vytvořit některé soubory:

<i>Soubor</i>	<i>Skript</i>	<i>Popis</i>	<i>Příklad</i>
Tagged	CrtTagged.bat	Převede korpus (v souboru <i>input</i> ) z formátu <i>desam</i> na formát požadovaný <i>taggerem</i> . Upraví i konce řádků z /CR/LF na /LF.	
Tagged1 Tagged2	Split.bat	Náhodně rozdělí soubor <i>tagged</i> na <i>tagged1</i> a <i>tagged2</i> (po řádcích).	
Utagged Utagged1 Utagged2	CrtUntagged.bat	Odstraní všechny značky v souborech <i>tagged</i> , <i>tagged1</i> a <i>tagged2</i> . Výsledek je v souborech s prefixem <i>U</i> .	
BWL	CrtBwl.bat	Ze souboru <i>Utagged</i> vytvoří seznam všech slov a seřadí je podle četnosti.	se , . v na
SWL	CrtSwl.bat	Ze souboru <i>Tagged1</i> (první polovina korpusu) vytvoří seznam, obsahující slovo a jeho nejčastější značku (a počet, pouze pro informaci).	, kE 18430 se k3 16862 . kE 15686 je k5 10256
<u>BGL</u>	CrtBGL.bat	Ze souboru <i>Utagged</i> vytvoří seznam slovních dvojic. Tento soubor je výstupní, používá se při samotném značkování.	výfukovými plyny své podpisy zeslabení negativního stoprocentní jistotou

Všechny tyto skripty lze dohromady vyvolat pomocí [learnUnknownInit.bat](#).

Samotné učení se spustí souborem [learnUnknown.bat](#). Vznikne tak soubor LR.

### **Učení kontextových pravidel**

Učení kontextových pravidel musí být spuštěno až po učení modulu pro neznámé slova, protože využívá některé soubory vzniklé při učení tohoto modulu.

<i>Soubor</i>	<i>Skript</i>	<i>Popis</i>	<i>Příklad</i>
TL	CrtTL.bat	Ze vstupního korpusu <i>tagged1</i> (první polovina) vytvoří seznam, obsahující slovo a seznam všech jeho možných značek. První uvedená značka je ta nejčastější, pořadí ostatních není určeno.	Jana k1gM k1gF Hradci k1gI podobným k2gI k2gN antika k1gF
FL	CrtFl.bat	Stejně jako <i>TL</i> , ale pro kompletní korpus <i>tagged</i> . Tento soubor je výstupní, používá se při samotném značkování.	
DTC	CrtDtc.bat	Aplikuje <i>tagger</i> , na soubor <i>Utagged2</i> . Vznikne tak korpus pro ověřování při učení (testovací množina).	

Všechny tyto skripty lze dohromady vyvolat pomocí `learnContextInit.bat`.

Samotné učení se spustí souborem `learnContext.bat`. Vznikne tak soubor `CLR`.

Všechny (čtyři) výstupní soubory, které jsou třeba ke značkování, je možné zkopírovat do zvoleného adresáře pomocí `copyOutput.bat`. Dočasné soubory vzniklé při učení maže skript `clearTmp.bat`.

## Kompilace modulů v C pod Win32

K překladu bylo použito GNU *Dev-C++* od [Bloodshed software](#), knihovna *Libgw32c* (kvůli funkci *getopt*). *Dev-c++* nepodporuje více projektů v jednom řešení (*solutions*), proto je pro každý projekt vytvořen samostatný adresář. Pro použití *taggeru* s českými značkami a pro úspěšné přeložení bylo třeba udělat následující změny v původním kódu:

### Start-state-tagger

- Ošetřen nesprávný počet parametrů (*null-exception*).
- Změněn algoritmus na počáteční značkování neznámých slov. Původně byla všechna slova začínající malým písmenem označena tagem *NN* a začínající velkým *NNP*. Nyní se všechna slova nejdříve označí tagem *kI*. (Obdobně byl upraven i skript *unknown-lexical-learn.prl*)

### Final-state-tagger

- Ošetřen nesprávný počet parametrů (*null-exception*).

### Tagger

- Přidán hlavičkový soubor *getopt.h* a knihovny *libgw32c.a*, *libole32.a*, *libuuid.a*. Kvůli funkci *getopt*, která není součástí *Dev-C++*. Tyto knihovny jsou distribuovány pod GNU licenci.
- Přidána přípona *.exe* k názvům modulů *start-state-tagger* a *final-state-tagger*.

### Contextual-rule-learn

- Ošetřen nesprávný počet parametrů (*null-exception*).
- Funkce *tmpnam* nahrazena, protože vytvářela názvy začínající lomítkem a programy z *cygwinu* hledají takovéto soubory v aktuálním adresáři. Pomocné soubory se jmenují např. *tmp1*, *tmp2*.
- Upravena cesta k programu *rm* z absolutní (*/bin/rm*) na relativní (v aktuálním adresáři).

## Dosažená přesnost

<i>Slovo</i>	<i>Přesnost (%), získaná přiřazením nejčastější značky</i>	<i>Přesnost taggeru</i>
Libovolné		95.06
Jana, vedoucí, je, se	90.097	95.6
Jana	74.5	89.775
vedoucí	39.131	60.624
je	89.851	95.843
se	90.541	95.595

Slovo *vedoucí* je v korpusu nejméně zastoupené (cca 50 vět) a navíc se klasifikuje do čtyř skupin. Proto má nejmenší úspěšnost.

Úspěšnost *taggeru* pro libovolné slovo se zdá být větší, než pro některá zaměřená slova. To je nejspíše způsobeno tím, že se mezi libovolné slova započítávají i oddělovače slov (tečka, čárka), uvozovky a podobně, které mají vysoký výskyt (viz *BWL*) a 100% úspěšnost.

Tento dokument je ve formě [html](#) [pdf](#).