

Součástí úkolu jsou Python3 skripty `assignment3_1.py` a `assignment3_2.py`.

## Task 1

Použili sme Brillův tagger z knihovny NLTK. Ako počiatkový tagger sme použili jednoduchý unigramový tagger, ktorý si zapamätá najčastejšie značku pre každé slovo tréningových dát. Použili sme počiatkovú množinu 18 pravidiel, rovnaké ako v ukážke v knihe<sup>1</sup>. Počet pravidiel je maximálne 200 (prednastavená hodnota NLTK).

## Výsledky

Angličtina	Přesnost	Čeština	Přesnost
1. stupeň	0.8878314603811958	1. stupeň	0.7611722195653329
2. stupeň	0.8877681340735528	2. stupeň	0.7784951378925554
3. stupeň	0.8591886684372233	3. stupeň	0.7446763540290621
4. stupeň	0.8916271557338613	4. stupeň	0.7794976095512296
5. stupeň	0.8927185098345708	5. stupeň	0.7772650116698494
Průměr	0.8838267856920804	Průměr	0.76822126654160594
$\sigma$	0.01247811683956135	$\sigma$	0.01355396354317235

Tabulka 1: Přesnost (accuracy) pro Brillův tagger pro oba jazyky.

## Task 2

### Nesupervizovaný HMM — Baum-Welch

Pre túto časť sme sa rozhodli použiť iba bigramové a unigramové štatistiky, pretože ak by sme nechceli použiť nuly “natvrdo” zabralo by to prílišne veľa času pre angličtinu, pre češtinu by to bolo jednoducho nerealizovateľné, časovo ani pamäťovo.

---

<sup>1</sup>Volne dostupné na <http://www.nltk.org/book/>