Министерство науки и высшего образования Российской Федерации

федеральное государственное автономное образовательное учреждение высшего образования «Национальный исследовательский университет ИТМО»

Факультет инфокоммуникационных технологий

Практическая работа №3 «Многоязычные корпуса и корпуса второго языка» по дисциплине: «Компьютерная лингвистика»

Выполнили:

студентки II курса ИКТ группы <u>К3242</u> Ф.И.О. <u>Боброва Мария Иосифовна, Шикалова Софья Сергеевна, Шечелова Дарья Александровна, Третьякова Екатерина Сергеевна</u>

Проверил: доцент ФТМИ Коцюба Игорь Юрьевич

ПРАКТИЧЕСКАЯ РАБОТА 3

Цель работы: проанализировать корпуса для одних из самых распространенных языков в мире, создать свой корпус выбранного языка

Задачи:

- 1. Найти корпуса для испанского, хинди, арабского, бенгальского, португальского языков
- 2. Выбрать один понравившийся язык
- 3. Изучить какие корпуса есть у этого языка
- 4. Составить свой корпус

Выполнение работы (Задание 7)

Одни из самых распространенных языков в мире - испанский (450 млн говорящих), хинди (больше 360 млн говорящих), арабский (больше 320 млн говорящих), бенгальский (больше 200 млн говорящих), португальский (ок. 200 млн говорящих).

Воспользовавшись каталогами CLARIN (www.clarin.eu) и ELRA (www.elra.info), поисковыми системами Гугл или Яндекс и Википедией, были найдены данные корпуса для каждого языка:

- 1. Корпус испанского языка
 - a) Ссылки представлены на сайте Национального корпуса русского языка. https://www.corpusdelespanol.org/ (корпус Марка Дэвиса) https://corp.hum.sdu.dk/cqp.es.html (поиск корпусов CorpusEye).
 - b) Два корпуса испанского языка, созданные Испанской королевской академией: современный корпус CREA (https://corpus.rae.es/creanet.html) и исторический корпус CORDE (https://corpus.rae.es/cordenet.html)

2. Корпус хинди

На сайте Национального корпуса русского языка представлена данная ссылка: https://www.cfilt.iitb.ac.in/~corpus/hindi/, однако доступ запрещен

3. Корпус арабского языка

На сайте Национального корпуса русского языка представлена данная ссылка: https://arabicorpus.byu.edu/

5. Корпус бенгальского языка

Разговорный корпус, представленный в виде аудиофайлов и отдельный документ с транскрибацией файлов: https://www.openslr.org/37

5. Португальский язык

На сайте Национального корпуса русского языка представлены данные ссылки:

https://www.corpusdoportugues.org/ (корпус Марка Дэвиса) https://corp.hum.sdu.dk/cqp.pt.html (сайт CorpusEye)

Выбранный нами язык для анализа и разработки корпуса – бенгальский.

Первым этапом был поиск датасета на бенгальском языке. Нами был найден архив, состоящий из 1891 аудиозаписи фраз на бенгальском языке на коммерчески направленные темы и tsv файл с транскрибированными аудиозаписями, выглядящий следующим образом:

```
এইচআর টেক্সটাইল বাংলাদেশের ভেতরে একাধিক আউটলেটের মাধ্যমে শাড়ি বাচ্চাদের পোশাব স্ট্যান্ডার্ড ব্যাংক এ ইসলামী ব্যাংকিং এর সুবিধা রমেছে লাফার্জ সুরমা সিমেন্ট সর্বাধিক ব্যবহৃত সিমেন্ট উৎপাদন করে পিপলস ইস্যারেস অব চায়না ছেষট্টি বছর আগে ব্যবসা চালু করে বয়গের অব চায়না ছেষট্টি বছর আগে ব্যবসা চালু করে বয়গেস একটি ইন্ডাম্বিয়াল গ্রুপ কেয়া ডেভেলপারস দেশের বিভিন্ন স্থানে স্থাপনা তৈরি করে থাকে ভেরাইজন কমিউনিকেশনস একটি আমেরিকান বড়ব্যান্ড ৪ টেলিযোগাযোগ কোম্পানি মোইল বিশ্বের ফ্রুত বিস্যুত্ত ইস্যারেস কোম্পানি সাইক পঞ্জিরারটিক প্রাইভেট কোম্পানি বিশ্বের কি প্রিরহন খাতে আইজটি কোম্পানি বিশ্বর করে প্রাইজি জারোমারিইলস একটি আতি পরিচিত নাম চায়না কমিউনিকেশনস ক্রমট্রকশন চানের বৃহত্তম বন্দর নির্মাণ ও নকশা এবং ড্রেজিং কোম্পানি পরিবহন খাতে আফতাব অটোমোবাইলস একটি আতি পরিচিত নাম চায়না কমিউনিকেশনস ক্রমট্রকশন চানের বৃহত্তম বন্দর নির্মাণ ও নকশা এবং ড্রেজিং কোম্পানি যোমা কনডেসড মিল্ক ইন্ডাস্ট্রিজ এর অতি জনপ্রিয় ব্রাহ্ত হলো ফ্রেশ মারবেনি তার কোম্পানির বিশ্বাস হিসেবে সততা উদ্ভাবন ও ঐক্য ধারণ করে ড্রেজিং কোম্পানি বিশ্বরাকান কমিকাল কাউসিলের সদস্যা লিবাটি মিউচুয়াল ইস্যারেস্ঠা কর্ম একটি শীর্ষস্থানীয় বেসরকারী সাধারণ বীমা কোম্পানী শাহজালাল ইসলামী ব্যাংক শিক্ষার্থীয় বেসরকারী সাধারণ বীমা কোম্পানী শাহজালাল ইসলামী ব্যাংক শিক্ষার্থীদের জন্য বৃত্তির করেছে জেবিএস একটি রাজিলিয়ান কেশাম্পানি যারা মাংস প্রক্রিয়াকর্মণ ব্যবসা করে থাকে বেক্সারেক ক্রেয়ার নাৎসি যুদ্ধাপরাধে অংশগ্রহণ করার জন্য প্রথম বিশ্বযুদ্ধার পর এবং দ্বিতীর বিশ্বযুদ্ধার সলিতে বাংলাদেশ দীপিং কর্পোরেশন একটি রাষ্ট্রীয় মালিকানাধীন ও ব্যবসা প্রতিষ্ঠান প্রসার করে আসতে বাংলাদেশ দিপিং কর্পোরেশন একটি রাষ্ট্রীয় মালিকানাধীন ও প্রবস্থা ক্রমারি কর্পোরেশন এব প্রথম ভি গ্রহ্মার সিমিটার বাংলাদেশ নির্মার সিমেন্ট মিনার বৃহত্তম ফার্মারিটিটিকাল কোম্পানি স্থানিক প্রবাস্থানিক এর শাখার সংখ্যা বিশি রয়েছে স্থানা কায়ে গ্রাংটিক এর শাখার সংখ্যা বিশি রয়েছে স্বিদ্যার সাংস্টিটিক নির ত্যানিজনিক প্রদান করে স্বানাদেশনি বিশ্বয়ার সাংস্টিটিক কলেন বেদেক প্রবাদিনিত এর শামার সংখ্যা বিশি রয়েছে স্থানার সাংস্কার বিদ্যানিক এর শামার সংখ্যা বিশি রয়েছে স্বান্ধান সিম্বার সাংস্টি বিদ্যার সাংস্টিক বিদ্যানিক এর শামার স্বান্ধানিক এ
   ban 00737 00012222450
 ban_00737_00015581920
ban_00737_00028634754
 ban_00737_00035050432
ban_00737_00068052117
ban_00737_00107291991
 ban 00737 00112921837
ban_00737_00120125731
ban_00737_00120232454
ban_00737_00142194715
 ban 00737 00176672702
ban 00737 00180288955
 ban_00737_00183428707
 ban 00737 00215879121
 ban 00737 00272549577
ban_00737_00295500853
ban_00737_00317863452
ban 00737 00359383339
ban_00737_00371172031
ban_00737_00376983060
 ban_00737_00399904268
ban_00737_00402143113
ban_00737_00422977360
ban 00737 00460641213
 ban 00737 00530142458
 ban 00737 00535389845
 ban 00737 00604373891
 ban_00737_00606350601
 ban 00737 00636200661
                                                                                                       দেশের অন্য যে কোন বেসরকারি ব্যাংকের থেকে পুরালী ব্যাংক এর শাখার সংখ্যা বেশি রয়েছে
সিব্লুফুর্মি গণপুঞ্জাতনী নিবের তাতীয় বহুত্ব জ্ঞাতীয় তেলু ক্যোম্পারি
 ban_00737_00669891346
```

Следующим этапом была подготовка данных для дальнейшего создания корпуса. Выбранным решением было преобразование tsv файла в txt формат. После преобразования файл был очищен от кодовых номеров аудиозаписей и в качестве разделителей между предложениями был выбран знак ";".

Файл приобрел следующий вид и был готов к дальнейшей обработке:

```
data bengali (1) – Блокнот
Файл Правка Формат Вид Справка
এইচাআর টেক্সটিইল বাংলাদেশের ভেতরে একাধিক আউটলেটের মাধ্যমে শাড়ি বাচ্চাদের পোশাক মহিলাদের পোশাক এবং অন্যান্য টেক্সটাইল পণ্য উৎপাদন ও বিপণন করে
স্ট্যান্ডার্ড ব্যাংক এ ইসলামী ব্যাংকিং এর সুবিধা রয়েছে ;
 লাফার্জ সুরমা সিমেন্ট সর্বাধিক ব্যবহৃত সিমেন্ট উৎপা<sup>ন</sup>
পিপলস ইন্স্যুরেন্স অব চায়না ছেষট্টি বছর আগে ব্যবসা চালু করে ;
বয়াগের একটি ইন্সার্সিয়াল গুল
কেয়া ডেভেলপারস দেশের বিভিন্ন স্থানে স্থাপনা তৈরি করে থাকে
ভেরাইজন কমিউনিকেশনস একটি আমেরিকান ব্রডব্যান্ড ও টেলিযোগাযোগ কোম্পানি
মেটলাইফ বিশ্বের দ্রুত বিস্তৃত ইন্স্যুরেন্স কোম্পানি
সাইফ পাওয়ারটেক প্রাইভেট কোম্পানি হিসেবে নিবন্ধিত এবং পরিচালিত
অলিম্পিক ইন্ডাস্ট্রিজ জীবনকে সহজ করে তুলেছে ;
বেইজিং অটোমোটিভ গ্রুপ একটি চীনা রাষ্ট্রীয় উদ্যোগ ও নিয়ন্ত্রণকারী কোম্পানি ;
পরিবহন খাতে আফতাব অটোমোবাইলস একটি অতি পরিচিত নাম ;
চায়না কমিউনিকেশনস কলটাকশন চীনের বহুতম বন্দর নির্মাণ ও নকশা এবং ডেজিং কোম্পানি
মেঘনা কনডেন্সড মিদ্ধ ইন্ডাস্ট্রিজ এর অতি জনপ্রির ব্র্যান্ড হলো ফ্রেশ
মারুবেনি তার কোম্পানির বিশ্বাস হিসেবে সততা উদ্ভাবন ও ঐক্য ধারণ করে ;
ডাউ কেমিক্যাল আমেরিকান কেমিক্যাল কাউন্সিলের সদস্য
লিবাটি মিউচুয়াল ইন্স্যুরেন্স গ্রুপ একটি আমেরিকান বিমা কোম্পনি
ফেডারেল ইন্যারেন্স বাংলাদেশের একটি শীর্ষস্থানীয় বেসরকারী সাধারণ বীমা কোম্পানী
 ণাহজালাল ইসলামী ব্যাংক শিক্ষার্থীদের জন্য বৃত্তির করেছে
জেবিএস একটি ব্রাজিলিয়ান কোম্পানি যারা মাংস প্রক্রিয়াকরণ ব্যবসা করে থাকে :
ঢাকা ব্যাংক ক্রেতার চাহিদার কথা মাথায় রেখে প্রকল্প নির্বাচন ও বাস্তবায়ন করে থাকে
বেক্সিয়কো টেক্সটাইল এব বানিজ্ঞািক প্রধান শাখা ধানমন্ট্রিতে অবস্থিত
বোরার নাথেদি যুদ্ধাপরাধে অংশপ্রথণ করার জন্য প্রথম বিশ্বযুদ্ধের পর এবং দ্বিতীয় বিশ্বযুদ্ধের সময় যুক্তরাট্রে ব্যবসা হারিয়েছিল
লিভে বাংলাদেশ দীর্ঘ সময় ধরে বাংলাদেশে অপারেশন ও বারসা প্রতিষ্ঠান প্রসার করে আসছে ;
বাংলাদেশ শিপিং কর্পোরেশন একটি রাষ্ট্রীয় মালিকানাখীন ও পরিচালিত সরকারি কর্পোরেশন এবং বাংলাদেশের সর্ববৃহৎ জাহাজ মালিক
ও এম ভি গ্রুপ একটি সমন্থিত আন্তর্জাতিক তেল ও গ্যাস কোম্পানি
সাইনোফার্ম চীনের বৃহত্তম ফার্মাসিউটিকাল কোম্পানি ;
প্রিমিয়ার সিমেন্ট মিলস এখন অধিক বিক্রিত সিমেন্ট ;
এশিয়া প্যাসিফিক জেনারেল ইন্সারেন্স গ্রাহকদের সর্বোচ্চ নিরাপন্তা প্রদান করে সর্বনিম্ন খরচ ও দাবির দ্রুত নিম্পন্তির মাধ্যমে
দেশের জন্য যে কোন বেসরকারি ব্যাংকের থেকে পুবালী ব্যাংক এর শাখার সংখ্যা বেশি রয়েছে ;
সিএনঙঙসি গণপ্রজাতন্ত্রী চীনের তৃতীয় বৃহত্তম জাতীয় তেল কোম্পানি
জিবিবি পাওয়ার বেসরকারি ইলেকট্রিসিটি জেনারেশন কোম্পানি
সাইনোকেয় এব মূল ব্যবসাব শক্তি কৃষি কেয়িকালে বিষেল এসৌট এবং আর্থিক সেবায
বিক্রসজাক্ষর মার্টকার নাংলাদেশের প্রথম আইক্রসড সন্দল্লান্ত ইম্পাভ কোম্পানি ;
এমঞ্জন্তভঞ্জিত ইন্যুরেন্স একটি জাপানি বীমা কোম্পানি ;
এবিবি বহুজন পরিচিত একটি বহুজাতিক কর্পোরেশন ;
এমআই সিমেন্ট ফ্যাক্ট্রীর ক্রাউন সিমেন্ট বাংলাদেশের সর্বাধিক রপ্তানিকৃত সিমেন্ট
ইউনিলিভার এর পণ্য বৈচিত্র্যে এক ও অপ্রতিদ্বন্দ্বী
ডয়চে ব্যাংক জার্মান মালিকানাধীন বৃহৎ ব্যাংক
বেক্সিমকো নিটিং উচ্চমানের বুনন বন্ধ উৎপাদন এবং বিক্রি করে ;
বাংলাদেশ সাবমেরিন ক্যাবল কোম্পানি এর লক্ষ্য হলো প্রধান কোম্পানি হিসেবে দেশব্যাপী সাবমেরিন ক্যাবলের টেলিযোগাযোগ সেবা প্রদান করা
বাজার মূলধন ও গ্রাহকদের পরিপ্রেক্ষিতে অস্ট্রেলিয়ায় চার বৃহত্তম আর্থিক প্রতিষ্ঠানের একটি ন্যাশনাল অস্ট্রেলিয়া ব্যাংক :
```

После этого в качестве инструмента для работы с данными был выбран bnlp toolkit¹. Этот инструмент был специально разработан для работы с бенгальским языком и реализован в качестве библиотеки python. Нами было выбрано провести исследование какие части речи встречаются чаще в данном датасете. Для реализации данного исследования было решено использовать метод Pos Tagging.

Возникшая проблема заключалась в том, что общепринятый список тэгов для английского языка кардинально отличается от тэгов для бенгальского, a bnlp toolkit практически не обладает никакой документацией и списка тэгов там не было.

Решение этой проблемы было найдено благодаря статье², в которой представлена таблица тэгов для бенгальского языка.

¹ Bengali Natural Language Processing(BNLP) — bnlp latest documentation

² ICO 2020 Paper 163.pdf

Сравнение POS Tags списков для английского и бенгальских языков:

Alphabetical list of part-of-speech tags used in the Penn Treebank Project:

Number	Tag	Description		
1.	CC	Coordinating conjunction		
2.	CD	Cardinal number		
3.	DT	Determiner		
4.	EX	Existential there		
5.	FW	Foreign word		
6.	IN	Preposition or subordinating conjunction		
7.	JJ	Adjective		
8.	JJR	Adjective, comparative		
9.	JJS	Adjective, superlative		
10.	LS	List item marker		
11.	MD	Modal		
12.	NN	Noun, singular or mass		
13.	NNS	Noun, plural		
14.	NNP	Proper noun, singular		
15.	NNPS	Proper noun, plural		
16.	PDT	Predeterminer		
17.	POS	Possessive ending		
18.	PRP	Personal pronoun		
19.	PRP\$	Possessive pronoun		
20.	RB	Adverb		
21.	RBR	Adverb, comparative		
22.	RBS	Adverb, superlative		
23.	RP	Particle		
24.	SYM	Symbol		
25.	TO	to		
26.	UH	Interjection		
27.	VB	Verb, base form		
28.	VBD	Verb, past tense		
29.	VBG	Verb, gerund or present participle		
30.	VBN	Verb, past participle		
31.	VBP	Verb, non-3rd person singular present		
32.	VBZ	Verb, 3rd person singular present		
33.	WDT	Wh-determiner		
34.	WP	Wh-pronoun		
35.	WP\$	Possessive wh-pronoun		
36.	WRB	Wh-adverb		

Рисунок 1 - POS Tags для английского языка

Table 1. Summary of POS tagsets

11 Tagset	Tagset Count	30 Tagset	Tagset Count
	44425	Common Noun (NC)	30819
Noun (N)		Proper Noun (NP)	7994
Noun (N)		Verbal Noun (NV)	2985
		Spatio-Temporal Noun (NST)	2627
Voul. (V)	14292	Main Verb (VM)	12062
Verb (V)	14292	Auxiliary Verb (VAUX)	2230
		Pronominals (PPR)	5137
	6409	Reflexive (PRF)	362
Pronoun (P)		Reciprocal (PRC)	15
	1	Relative (PRL)	448
	1	WH Pronoun (PWH)	447
N'1 M1'C (T)	14332	Adjective (JJ)	9377
Nominal Modifier (J)		Quantifier (JQ)	4955
		Absolute Demostrative (DAB)	2421
Demonstrative (D)	2876	Relative Demostrative (DRL)	400
` '		WH Demostrative (DWH)	55
A 1 1 (A)	3965	Adverb of Manner (AMN)	1995
Adverb (A)		Adverb of Location (ALC)	1970
Doubletele (I)	573	Verbal Participle (LV)	72
Participle (L)		Conditional Participle (LC)	501
Post Position (PP)	3989	Post Position (PP)	3989
, ,		Coordinating Particle (CCD)	2899
	6704	Subordinating Particle (CSB)	2051
Particle (C)		Classifier Particle (CCL)	324
, ,		Interjection (CIN)	59
		Others (CX)	1371
Punctuation (PU)	13519	Punctuation (PU)	13519
, , ,	4348	Foreign Word (RDF)	1873
Residual (R)		Symbol (RDS)	1968
` ′		Others (RDX)	507

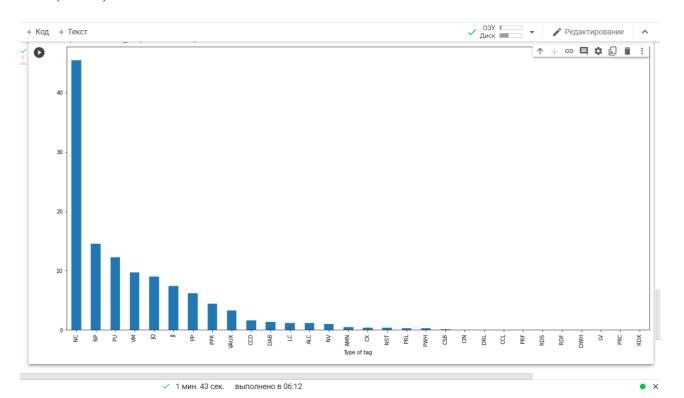
Pucyнок 2 - POS Tags для бенгальского языка

Расшифровка морфологических единиц для английского и бенгальского языка значительно различаются, например, есть более специфические расшифровки для бенгальского языка (в английском – adverb (RB), в бенгальском – 2 различных наречия (AMN, ALC). То же самое касается большинства частей речи.

Далее нами был написан код на языке Python по созданию корпуса бенгальского языка, который позволяет посчитать процент содержания каждого тэга в тексте. Ввести нужный тэг можно с клавиатуры. В качестве обучающий модели была взята модель для Pos Tagging, предложенная на github разработчиков bnlp toolkit. Данная модель позволяет достигнуть 80% точности.

```
# You will be a provided by the continue of th
```

Также была составлена столбчатая диаграмма, которая показывает что тремя наиболее используемыми тэгами являются NC(Common Noun), NP(Proper Noun) и PU(Punctuation). Помимо этого есть тэги, которые и вовсе не встречаются в корпусе – это RDS(Symbol), RDF(Foreign Word), DWH(WH Demonstrative), LV(Verbal Participle), PRC(Reciprocal), RDX(Others).



Выводы:

Малое количество таких тэгов как DWH, LV и PRC в целом характерно для датасетов на бенгальском языке (если полагаться на таблицу из статьи, где подсчитано количество тэгов для 7390 предложений на бенгальском, поэтому их отсутствие в выбранном нами датасете неудивительно), тогда как отсутствие тэгов как RDF, RDS и RDX можно считать спецификой коммерчески направленных предложений на бенгальском языке, а отсутствие тэга RDS легко объясняется тем, что предложения являются транскрибированными аудиозаписями.

Источники:

- 1. file:///C:/Users/work11pro1/Downloads/ICO_2020_Paper_163%20(1).pdf
- 2. https://catalog.ldc.upenn.edu/docs/LDC2010T24/Annotation_Guidelines_for_Hindi.pdf