

rus	300	17	29	22	2	3	29
bel	17	300	7	8	1	2	22
bulg	29	7	300	28	1	2	14
serb	22	8	28	300	1	2	11
tat	2	1	1	1	300	8	0
kaz	3	2	2	2	8	300	1
ukr	29	22	14	11	0	1	300
	rus	bel	bulg	serb	tat	kaz	ukr

rus	300	180	233	194	153	128	210
bel	180	300	159	144	125	123	187
bulg	233	159	300	209	138	119	198
serb	194	144	209	300	132	112	181
tat	153	125	138	132	300	164	128
kaz	128	123	119	112	164	300	120
ukr	210	187	198	181	128	120	300
	rus	bel	bulg	serb	tat	kaz	ukr

rus	300	65	168	119	38	22	120
bel	65	300	53	52	31	26	83
bulg	168	53	300	122	33	22	92
serb	119	52	122	300	22	20	87
tat	38	31	33	22	300	64	32
kaz	22	26	22	20	64	300	23
ukr	120	83	92	87	32	23	300
	rus	bel	bulg	serb	tat	kaz	ukr

Какие языки оказались очень похожи друг на друга, а какие нет? Иными словами, какие языки проще различить по нграммам, а какие труднее? Интерпретируйте результат.

На основе полученных данных можно сделать ожидаемые выводы, что русский очень схож с украинским, болгарским и сербским языками, а казахский и татарский языки, напротив, ощутимо отличаются от языков славянской группы.

Особенно это сходство заметно при анализе биграмм: у болгарского и русского языков 233 общих высокочастотных биграмм из 300, у украинского с русским - 210, у сербского и болгарского - 209. Это очень высокие показатели, которые указывают на высокую степень схожести языков, что сильно затруднит определение языка по биграммам. Количество общих высокочастотных биграмм у казахского и татарского языков значительно ниже - в диапазоне от 112 - 164, однако это тоже довольно высокие показатели (почти половина от общего количества биграмм). Причем между собой татарский и казахский довольно сильно похожи.

В случае с триграммами картина выглядит несколько лучше, однако языки, демонстрировавшие высокую степень сходства между собой на основе анализа биграмм, в данном случае также обладают высокими показателями: на первом месте количество общих триграмм у русского и болгарского языка - 168, также сильно похожи болгарский и сербский - 122 общих триграмма. Точно также казахский и татарский языки на фоне общей картины имеют меньшую степень схожести с другими языками, и при этом демонстрируют высокую степень сходства между собой.

Анализ частотных слов выглядит наиболее оптимальным для определения языка, так как общих высокочастотных слов у этих языков значительно меньше, чем общих биграмм и триграмм: максимальный показатель - 29.

При этом матрица схожести языков еще раз подтверждает выводы, сделанные на основе анализа общих биграмм и триграмм: больше всего похожи между собой русский и болгарский (возможно, это объясняется практически одинаковым алфавитом в этих двух языках), сильное сходство русского, сербского, украинского между собой, и ощутимая дистанция татарского и казахского языков.

Выводы:

- анализ трех матриц показал, что языки, принадлежащие славянской группе, ожидаемо демонстрируют высокую степень схожести;
- на основе сравнения трех матриц между собой можно сделать вывод о том, что менее всего целесообразно определять язык по биграммам, а наиболее оптимально - по отдельным высокочастотным словам.