

Gramáticas Locais para Reconhecimento de Construções com Verbo Suporte em Português

Luís Enrique Santos Prado Vereau, Juliana Pinheiro Campos Pirovani

¹Departamento de Computação - Universidade Federal do Espírito Santo;
Alegre, ES - Brasil

luis.vereau@edu.ufes.br, juliana.campos@ufes.br

Abstract. *Natural Language Processing is an interdisciplinary subarea of Computer Science and Linguistics that aims to study the generation, representation and understanding of natural language by computers. This article describes the automatic generation of Local Grammars (LGs) for Supporting Verb Constructions (SVC) from six Lexicon-Grammar tables that describe a total of 468 CVS. For this, six parameterized graphs were created using the tool Unitex to extract information from these tables. The LGs generation was done through shell scripts and Unitex. The generated LGs were applied at the corpus aTribuna, where 211 structures in the form of the searched SVC were found.*

Resumo. *O Processamento de Linguagem Natural é uma subárea interdisciplinar da Ciência da Computação e Linguística que tem como objetivo o estudo da geração, representação e compreensão da linguagem natural por computadores. Este artigo descreve a geração semiautomática de Gramáticas Locais (GLs) para Construções com Verbo Suporte (CVS) a partir de seis tábuas do Léxico-Gramática que descrevem um total de 468 CVS. Para isto, foram criados seis grafos parametrizados utilizando a ferramenta Unitex para a extração de informações destas tábuas. A geração das GLs foi feita por meio de shell scripts e do Unitex. As GLs geradas foram aplicadas no corpus aTribuna, onde foram encontradas 211 estruturas no formato das CVS buscadas.*

1. Introdução

O Processamento de Linguagem Natural (PLN) é uma subárea interdisciplinar da Ciência da Computação e Linguística que se dedica a geração, representação e compreensão de linguagem natural de forma automática. Os próprios linguistas podem se beneficiar do PLN por meio das ferramentas construídas pelos profissionais da Computação, da mesma forma em que, simetricamente, a qualidade do PLN pode depender da descrição da língua pelos linguistas [Picoli et al. 2015]. [Chowdhury 2003] descreve PLN como sendo “uma área de pesquisa e aplicação que explora como computadores podem ser usados para entender e manipular texto e fala em linguagem natural para fazer coisas úteis”.

Devido a necessidade atual de compartilhamento e compreensão de informação na era da internet, essa área tem ganhado destaque cada vez maior no meio acadêmico e na indústria. A grande quantidade de informação disponível atualmente em textos de escrita livre (não estruturados) precisa ser tratada para uso em aplicações que buscam informações específicas a partir deles. Exemplos de aplicações relevantes que necessitam das técnicas de PLN são sistemas de perguntas e respostas, tradução automática e reconhecimento de entidades nomeadas.

Entre os diversos problemas encontrados pelas máquinas ao processar a linguagem natural, há a interpretação de expressões com significado não composicional, que não podem ter seu significado extraído analisando o sentido literal de suas palavras individualmente. Expressões Cristalizadas (EC) como “João comprou no mercado negro” e Construções com Verbo Suporte (CVS) como “João tem sangue frio para blefar” são exemplos dessas expressões. Este trabalho tem CVS como objeto de estudo, [Picoli 2020] expõe as diferenças entre EC e CVS.

CVS podem ser definidas como expressões compostas por um verbo que atua como Verbo Suporte (Vsup) e uma unidade predicativa não-verbal que pode ser um nome predicativo (Npred) como em “ter lábia”(Vsup+Npred), um adjetivo (adj) como na construção “estar liso”(Vsup+adj), ou uma expressão que se comporta como adjetivo (Expadj), “estar azul de fome”(Vsup+Expadj), por exemplo [Picoli 2020].

[Flores 2020] analisou ocorrências do verbo *dar* em duas coleções de textos em português brasileiro, o *corpus* do projeto Fala Goiânia¹ e o *corpus* do Grupo Discurso & Gramática². Notou-se que das 190 ocorrências do verbo no *corpus* do projeto Fala Goiânia, em 103 (54,21%) ele atuava como Vsup, e no *corpus* do Grupo Discurso & Gramática, das 70 ocorrências encontradas, em 47 (67,14%) o verbo assume papel de Vsup, indicando grande presença das CVS na língua portuguesa do Brasil.

Desta forma, o reconhecimento automático ou semiautomático de CVS é importante para o PLN pela riqueza de significado de seus elementos, a dificuldade de entendimento da não-composicionalidade pelo computador, bem como pelo frequente aparecimento de expressões desse tipo no português. A geração automática de resumos, tradução de máquina e *chatbots* [Tan et al. 2021] são possíveis tarefas nas quais o reconhecimento dessas expressões poderá ser aplicado.

Um dos métodos para descrever CVS é o Léxico-Gramática [Gross 1975] que consiste em criar tabelas, também chamadas de tábuas, que detalham um conjunto de expressões e possíveis variações para certos elementos, como verbo utilizado, ausência ou existência de negação, comparação, intensificação, dentre outras.

Uma forma de reconhecer expressões com características sintáticas e semânticas em comum, como as CVS, são as Gramáticas Locais (GLs). As Gramáticas Locais [Gross 1997] são “gramáticas de estados finitos ou autômatos de estados finitos que representam conjuntos de expressões de uma língua natural”.

Este trabalho tem como objetivo a extração de informações presentes nas seis tábuas do Léxico-Gramática que descrevem CVS em [Picoli 2020] para geração semi-automática de GLs que sejam capazes de reconhecer essas expressões e anotar seus significados. O trabalho busca, por fim, fornecer recursos que possam ser utilizados em aplicações de PLN e auxiliem tanto profissionais da computação quanto da linguística.

2. Metodologia

Unitex³ é um conjunto de software livres para PLN que permite, além do pré-processamento de textos, a construção de GLs e a construção automática de GLs a partir

¹<https://gef.letras.ufg.br/p/11947-projetos-tematicos>

²<https://discursoegramatica.wordpress.com/>

³<https://unitexgramlab.org/pt/>

de uma tabela (neste caso, tábua do Léxico-Gramática) e de um grafo parametrizado.

As GLs no Unitex são representadas por grafos como o apresentado na Figura 1 que reconhece a estrutura [Nome Próprio (reconhecido pelo código lexical <N+Pr >) + não (opcional) + verbo *ser*, *ter*, *permanecer* ou *continuar* (reconhecido pelo código lexical <Verbo.V >, como <ser.V >) + coração grande].

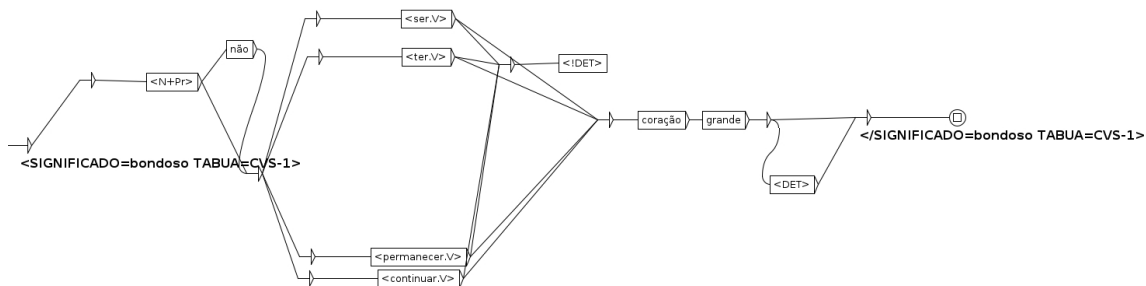


Figura 1. SubGL CVS-1

Um exemplo de tábua do Léxico-Gramática é apresentado na Figura 2. Na linha 3 desta tábua está descrita a CVS representada pela subGL da Figura 1.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325	326	327	328	329	330	331	332	333	334	335	336	337	338	339	340	341	342	343	344	345	346	347	348	349	350	351	352	353	354	355	356	357	358	359	360	361	362	363	364	365	366	367	368	369	370	371	372	373	374	375	376	377	378	379	380	381	382	383	384	385	386	387	388	389	390	391	392	393	394	395	396	397	398	399	400	401	402	403	404	405	406	407	408	409	410	411	412	413	414	415	416	417	418	419	420	421	422	423	424	425	426	427	428	429	430	431	432	433	434	435	436	437	438	439	440	441	442	443	444	445	446	447	448	449	450	451	452	453	454	455	456	457	458	459	460	461	462	463	464	465	466	467	468	469	470	471	472	473	474	475	476	477	478	479	480	481	482	483	484	485	486	487	488	489	490	491	492	493	494	495	496	497	498	499	500	501	502	503	504	505	506	507	508	509	510	511	512	513	514	515	516	517	518	519	520	521	522	523	524	525	526	527	528	529	530	531	532	533	534	535	536	537	538	539	540	541	542	543	544	545	546	547	548	549	550	551	552	553	554	555	556	557	558	559	560	561	562	563	564	565	566	567	568	569	570	571	572	573	574	575	576	577	578	579	580	581	582	583	584	585	586	587	588	589	590	591	592	593	594	595	596	597	598	599	600	601	602	603	604	605	606	607	608	609	610	611	612	613	614	615	616	617	618	619	620	621	622	623	624	625	626	627	628	629	630	631	632	633	634	635	636	637	638	639	640	641	642	643	644	645	646	647	648	649	650	651	652	653	654	655	656	657	658	659	660	661	662	663	664	665	666	667	668	669	670	671	672	673	674	675	676	677	678	679	680	681	682	683	684	685	686	687	688	689	690	691	692	693	694	695	696	697	698	699	700	701	702	703	704	705	706	707	708	709	710	711	712	713	714	715	716	717	718	719	720	721	722	723	724	725	726	727	728	729	730	731	732	733	734	735	736	737	738	739	740	741	742	743	744	745	746	747	748	749	750	751	752	753	754	755	756	757	758	759	760	761	762	763	764	765	766	767	768	769	770	771	772	773	774	775	776	777	778	779	780	781	782	783	784	785	786	787	788	789	790	791	792	793	794	795	796	797	798	799	800	801	802	803	804	805	806	807	808	809	810	811	812	813	814	815	816	817	818	819	820	821	822	823	824	825	826	827	828	829	830	831	832	833	834	835	836	837	838	839	840	841	842	843	844	845	846	847	848	849	850	851	852	853	854	855	856	857	858	859	860	861	862	863	864	865	866	867	868	869	870	871	872	873	874	875	876	877	878	879	880	881	882	883	884	885	886	887	888	889	890	891	892	893	894	895	896	897	898	899	900	901	902	903	904	905	906	907	908	909	910	911	912	913	914	915	916	917	918	919	920	921	922	923	924	925	926	927	928	929	930	931	932	933	934	935	936	937	938	939	940	941	942	943	944	945	946	947	948	949	950	951	952	953	954	955	956	957	958	959	960	961	962	963	964	965	966	967	968	969	970	971	972	973	974	975	976	977	978	979	980	981	982	983	984	985	986	987	988	989	990	991	992	993	994	995	996	997	998	999	1000
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325	326	327	328	329	330	331	332	333	334	335	336	337	338	339	340	341	342	343	344	345	346	347	348	349	350	351	352	353	354	355	356	357	358	359	360	361	362	363	364	365	366	367	368	369	370	371	372	373	374	375	376	377	378	379	380	381	382	383	384	385	386	387	388	389	390	391	392	393	394	395	396	397	398	399	400	401	402	403	404	405	406	407	408	409	410	411	412	413	414	415	416	417	418	419	420	421	422	423	424	425	426	427	428	429	430	431	432	433	434	435	436	437	438	439	440	441	442	443	444	445	446	447	448	449	450	451	452	453	454	455	456	457	458	459	460	461	462	463	464	465	466	467	468	469	470	471	472	473	474	475	476	477	478	479	480	481	482	483	484	485	486	487	488	489	490	491	492	493	494	495	496	497	498	499	500	501	502	503	504	505	506	507	508	509	510	511	512	513	514	515	516	517	518	519	520	521	522	523	524	525	526	527	528	529	530	531	532	533	534	535	536	537	538	539	540	541	542	543	544	545	546	547	548	549	550	551	552	553	554	555	556	557	558	559	560	561	562	563	564	565	566	567	568	569	570	571	572	573																																																																																																																																																																																																																																																																																																																																																																																																																																											

As GLs geradas pelos grafos parametrizados foram aplicadas no *corpus* composto por um conjunto de textos publicados pelo jornal do Espírito Santo aTribuna⁴. Esse *corpus* contém textos jornalísticos de gêneros variados, tendo sido utilizado em [Santiago 2022]. Também foram criados *shell scripts* que automatizaram o processo de criação e aplicação das GLs ao chamar programas do Unix responsáveis pela *tokenização* dos textos (Tokenize), aplicação de um Grafo Parametrizado em uma Tábua do Léxico-Gramática (Table2Grf), compilação das GLs (Grf2Fst2), dentre outros.

3. Resultados

Foram construídos 6 grafos parametrizados e geradas semiautomaticamente 473 GLs para reconhecimento de CVS. Para o *corpus* escolhido, foram encontradas 211 correspondências com estruturas no formato das CVS buscadas.

Algumas das expressões corretamente identificadas foram: "Vaz tem carta branca", "três é show de bola", "eu sou o máximo", "vida é um mar de rosas", "senão fica o dito pelo não dito", "público não esteve lá essas coisas", "Eduardo está entre a cruz e a espada", "Motta é do ramo" e "Sandra tem os dias contados", sendo "ter os dias contados" a CVS mais frequente, aparecendo 21 vezes (9,95%).

Observou-se também a presença de falso-positivos dentre as construções identificadas, i.e, expressões que foram identificadas como CVS, mas não são CVS. Por exemplo, anotou-se a expressão "o corpo foi achado", falso positivo para a CVS *ser um achado*, como em "o livro é um achado". Notou-se que a ocorrência deste falso-positivo em particular foi devido à descrição na tábua CVS-1cop que permite essas construções sem o artigo indefinido. Sabendo que este artigo é necessário para esta CVS, a tábua poderia ser alterada para corrigir essa situação melhorando a precisão da GL gerada.

Outra CVS que apresentou falso-positivos foi *ser o de menos*, que indica pouca relevância como em "essa informação será o de menos". Um dos falso-positivos identificados foi "A academia também observou que a taxa de acréscimo de novas reservas em todo o mundo era de menos", que indica que taxa de acréscimo de novas reservas em todo o mundo foi menor que certo valor, apresentando, desta forma, significado bem diferente ao da CVS.

Faz-se necessário, então, analisar os grafos parametrizados novamente a fim de mitigar a aparição de falso-positivos.

4. Conclusão

Neste trabalho foram geradas 473 GLs a partir de tábuas do Léxico-Gramática e grafos parametrizados construídos no Unix. As GLs identificaram no *corpus* aTribuna 211 estruturas correspondentes a CVS.

Como metas futuras, os grafos parametrizados construídos até então sofrerão novas análises e serão feitas melhorias. Igualmente serão feitas novas análises nas tábuas do Léxico-Gramática a fim de remover possíveis falhas. Também será realizada a análise de resultados, avaliando-se a precisão das correspondências encontradas.

Pode-se avaliar também em quais gêneros textuais houve mais ocorrências de CVS e analisar se, nesses textos, há alguma CVS não descrita nas tábuas.

⁴<https://tribunaonline.com.br/>

Referências

- Chowdhury, G. (2003). Natural language processing. *Annual review of information science and technology*, pages 51–89.
- Flores, E. (2020). A construção-suporte no português brasileiro.
- Gross, M. (1975). Méthodes en syntaxe: régime des constructions complétives.
- Gross, M. (1997). The construction of local grammars. *Finite-state language processing*, pages 329–354.
- Picoli, L. (2020). Contínuo e limite entre expressão cristalizada e construção com verbo-suporte à luz do léxico-gramática.
- Picoli, L. et al. (2015). Uso de uma ferramenta de processamento de linguagem natural como auxílio à coleta de exemplos para o estudo de propriedades sintático-semânticas de verbos. *Linguamática*, pages 35–44.
- Santiago, D. (2022). Gramáticas locais para reconhecimento de expressões cristalizadas em português. *Relatório Parcial de Pesquisa, Programa Institucional de Iniciação Científica 2021/2022*.
- Tan, K. et al. (2021). Review on light verb constructions in computational linguistics. pages 25–26.