

SUBTLEX_PT: Uma nova medida de frequência lexical para o português europeu baseada em legendas de filmes e séries televisivas

Ana Paula Soares¹, João Machado¹, Ana Costa¹, Alberto Simões², José João de Almeida², Álvaro Iriarte² & Montserrat Comesaña¹

{¹Escola Psicologia, ²Departamento de Informática, ³Instituto de Letras e Ciências Humanas}, Universidade do Minho

A validade concorrente das medidas de frequência lexical utilizadas em estudos psicolinguísticos tem sido recentemente questionada. De uma forma geral esses estudos têm demonstrado que as medidas de frequência e de diversidade contextual (número de filmes em que a palavra ocorre) extraídas a partir de legendas de filmes e séries televisivas explicam uma percentagem significativamente maior de variância da precisão e dos tempos de reconhecimento e nomeação de palavras do que outras medidas de frequência classicamente utilizadas. Neste trabalho apresentamos uma nova medida de frequência lexical e de diversidade contextual para ≈136.000 palavras do português europeu extraídas a partir de um *corpus* de ≈78 milhões de palavras derivados de ≈17.500 legendas de filmes e séries televisivas obtidos a partir do OPUS (<http://opus.lingfil.uu.se/>).

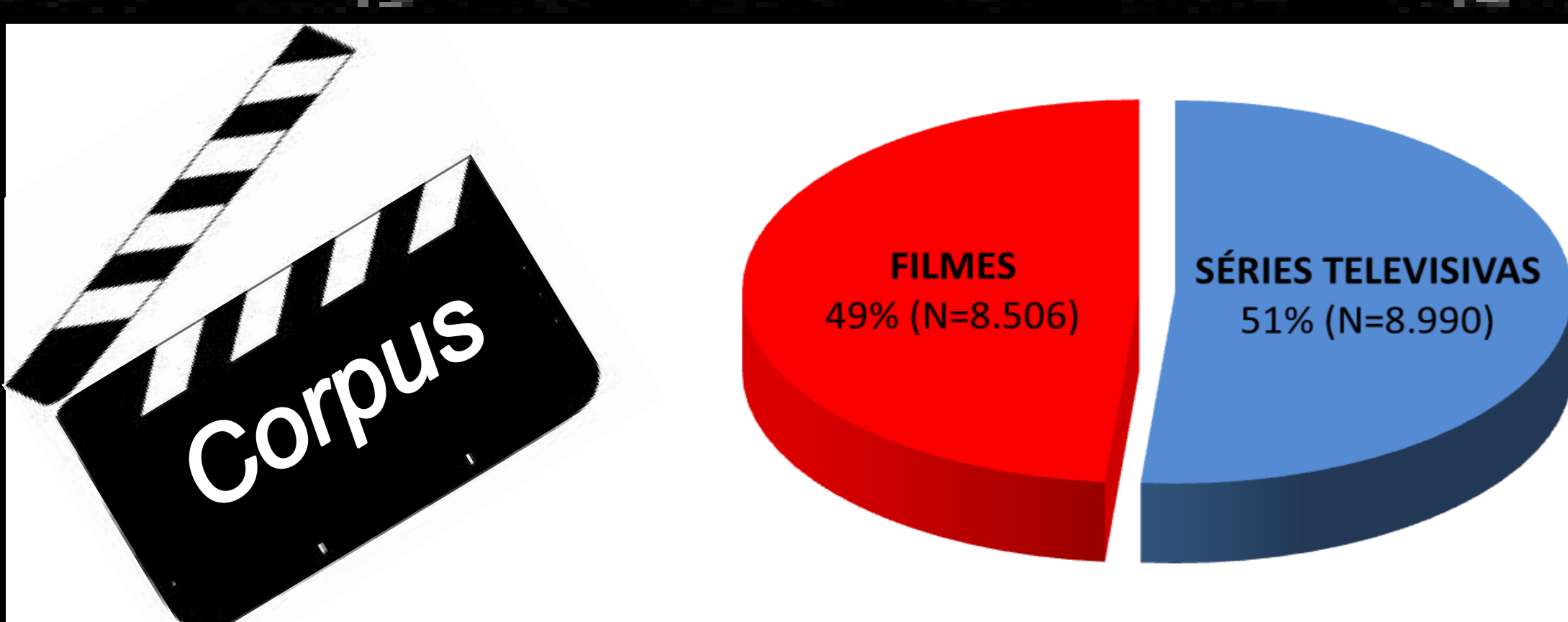


Figura 1. Distribuição das 17.496 legendas no *corpus* SUBTLEX_PT.

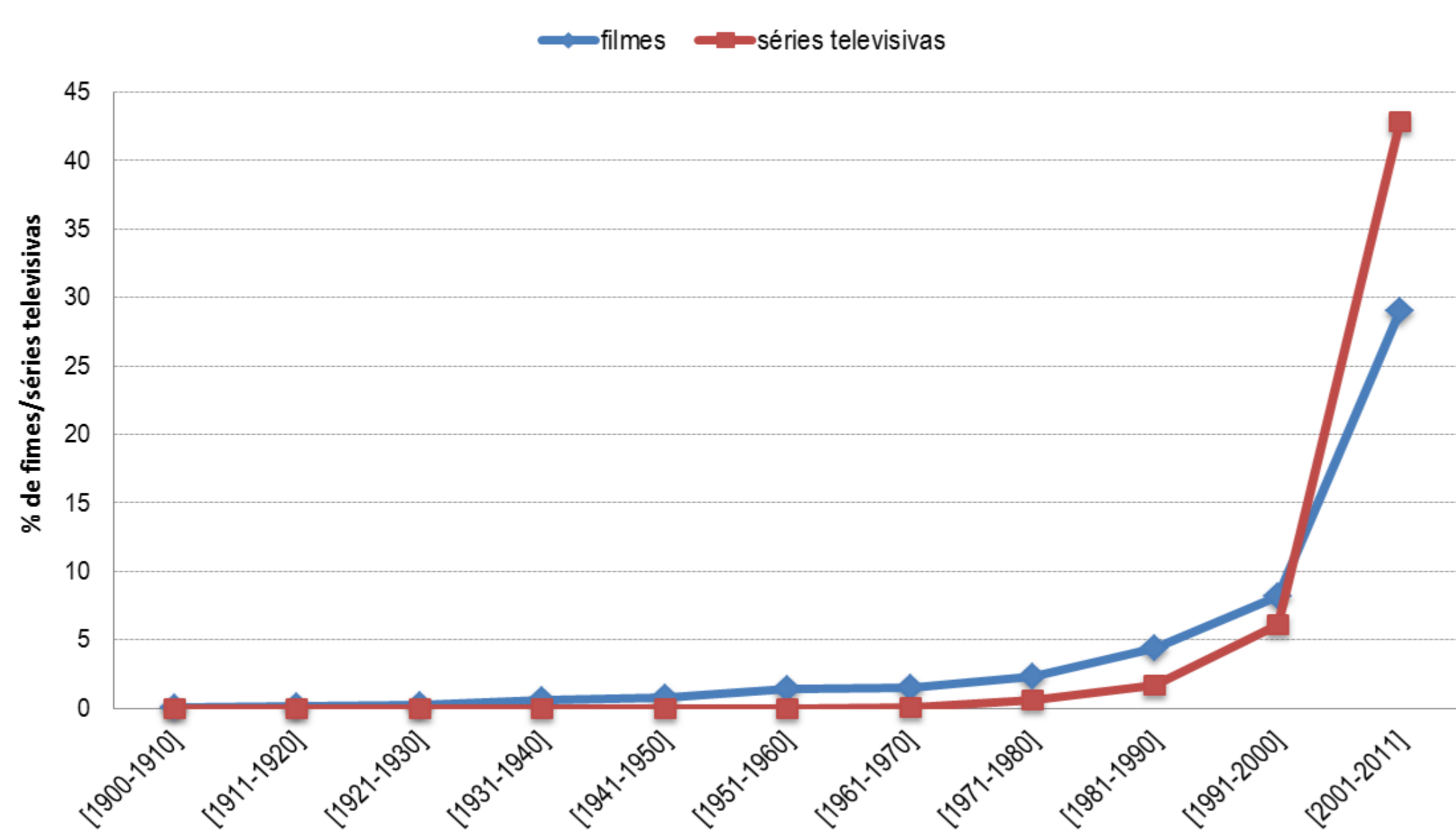


Figura 2. Distribuição dos filmes e séries televisivas do *corpus* SUBTLEX_PT por década.

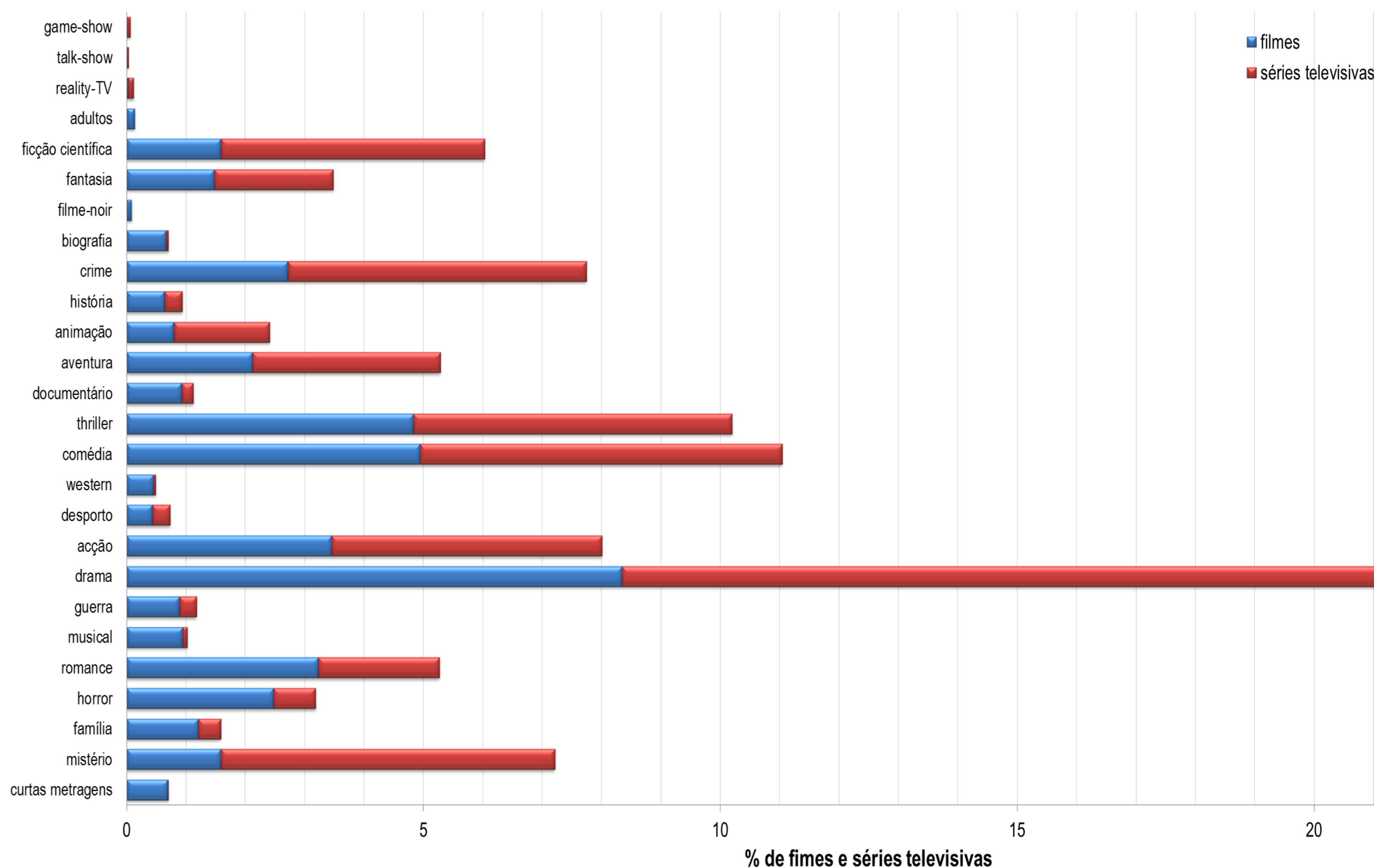


Figura 3. Distribuição dos filmes e séries televisivas do *corpus* SUBTLEX_PT por género cinematográfico (<http://www.imdb.com/>).



DC_{cont}: número de filmes e séries televisivas nos quais a palavra ocorre (i.e., num valor máximo de 17.496 filmes e séries televisivas).

FREQ_{milhão}: medida standard de frequência por milhão de palavras que toma em consideração a dimensão do *corpus*. É apresentada com 2 dígitos de precisão para não perder informação das contagens absolutas.

LOG10_{cont}: valor que resulta do cálculo do logaritmo de base 10 da FREQ_{cont}+1. Como a medida FREQ_{cont} se baseia num *corpus* de 78.402.091 palavras, um valor LOG10=0,3 corresponde a palavras que ocorrem apenas uma vez no *corpus* e LOG10>5 que ocorrem mais de 100.000 vezes no *corpus*. É apresentada com 4 dígitos de precisão.

PALAVRA: contém 135.598 formas flexionadas (flexões verbais e nominais) do português europeu que ocorrem no *corpus* SUBTLEX_PT. Da base lexical fazem parte todas as formas distintas, não se diferenciando os casos das homógrafas não homófonas (ex. "forma" ['fɔrma] e "forma" ['fɔrme]) e as palavras homónimas (ex. "além" [nome] e "além" [advérbio]), que constituem entrada única na base.

FREQ_{cont}: número de vezes que a palavra ocorre no *corpus* SUBTLEX_PT (i.e., no total de 78.402.091 palavras).

	A	B	C	D	E	F	G
	PALAVRA	FREQ _{cont}	DC _{cont}	FREQ _{milhão}	LOG10 _{cont}	DC _%	LOG10 _{DC}
1	a	2701655	9246	34458,97	6,4316	52,85	3,9660
2	de	2096154	9221	26735,95	6,3214	52,70	3,9648
3	o	2584427	9209	32963,75	6,4124	52,63	3,9643
4	que	2999219	9208	38254,32	6,4770	52,63	3,9642
5	no	551497	9195	7034,21	5,7415	52,55	3,9636
6	para	912398	9195	11637,42	5,9602	52,55	3,9636
7	se	838421	9194	10693,86	5,9235	52,55	3,9636
8	por	606118	9191	7730,89	5,7826	52,53	3,9634
9	como	438050	9176	5587,22	5,6415	52,45	3,9627
10	me	765240	9176	9760,45	5,8838	52,45	3,9627
11	nos	218977	9111	2793,00	5,3404	52,07	3,9596
12	ser	220595	9091	2813,64	5,3436	51,96	3,9587
13	nada	173218	9073	2209,35	5,2386	51,86	3,9578
14	está	668064	9067	8521,00	5,8248	51,82	3,9575
15	vamos	260254	9059	3319,48	5,4154	51,78	3,9571
16	ver	134103	9020	1710,45	5,1274	51,55	3,9553
17	da	425268	9010	5424,19	5,6287	51,50	3,9548
18	todos	130561	8976	1665,27	5,1158	51,30	3,9531
19	era	163585	8970	2086,49	5,2137	51,27	3,9528
20	nunca	122106	8964	1557,43	5,0867	51,23	3,9526
21	e	1299762	8961	16578,15	6,1139	51,22	3,9524
22	vez	109379	8959	1395,10	5,0389	51,21	3,9523
23	porque	188080	8916	2398,92	5,2743	50,96	3,9502
24	esta	140239	8900	1788,72	5,1469	50,87	3,9494
25	dos	128471	8888	1638,62	5,1088	50,80	3,9489
26	este	127787	8861	1629,89	5,1065	50,65	3,9475
27	ir	109193	8852	1392,73	5,0382	50,59	3,9471
28	os	476814	8769	6081,65	5,6783	50,12	3,9430
29	mas	542404	8766	5642,76	5,6458	50,10	3,9429
30	te	419723	8763	6628,94	5,7158	50,09	3,9427
31	com	600065	8723	7653,69	5,7782	49,86	3,9407

DC_%: valor que indica a percentagem de filmes e séries televisivas em que a palavra ocorre. É apresentada com 2 dígitos de precisão.

LOG10_{DC}: valor que resulta do cálculo do logaritmo de base 10 da DC_{cont}+1. Como a medida DC_{cont} se baseia em 17.496 legendas de filmes e séries televisivas, um valor LOG10_{DC}<1,08 corresponde a palavras que ocorrem em menos de 10 filmes e séries televisivas e LOG10>3 que ocorrem em mais de 1.000 filmes e séries televisivas. É apresentada com 4 dígitos de precisão.

para mais informações:
asoares@psi.uminho.pt