

An Inductive Classification of Italian Clusters



Matteo Pascoli
Università di Verona
Université de Toulouse - Le Mirail
matteo.pascoli@univr.it

words

Input data is a lexicon of Italian lemmas, with phonetic transcription made with heuristics on data from Italian online dictionaries (Treccani, Hoepli, Sabatini-Coletti) and refined by hand.

clusters

The words are divided around nuclei (vowels) to fill up a database of consonantic (including glides) clusters. Clusters are classified as initial, intervocalic, final.

sonority patterns

Each segment of the clusters is mapped to a sonority value. These sonority values are chained into a string, or sonority pattern. The clusters are then grouped according to their sonority pattern.

sonority contours

Each sonority pattern is linked with a sonority contour. Sonority contours are strings of symbols that represent the variation of sonority from a segment to the next. These symbols are: # for word boundary, / for rising sonority, \ for falling sonority, and - for plateau (same sonority level).

syllabification

I tried to use contours as indices for deciding where to break each cluster, that is, for the syllabification of the words. Most contours are sufficient indices, with the exception of contours 5 (V/), 8 (V//), and 16 (\\V). For clusters belonging to these contour groups, the sonority patterns are to be examined as indices for syllabification.

• syllable structure summary

structure	tokens	initial	median	final	isolate
CV	64014	7970 'dɛ.a	30559 e'ro.e	25457 'a.ja	28 ma
CVC	21936	7778 mu'la	14046 e'let.ta	83 'ɛ.ros	29 bar
CCV	8907	2594 'tri.o	5123 e'bre.o	1181 'ɔ.ljo	9 psi
VC	4881	4562 'is.sa	311 du'el.lo	6 ra'is	2 il
V	3582	1853 'ɔ.zo	1011 'ɔ.a.zi	713 're.o	5 o
CCVC	3488	2167 'skol.lo	1305 e'trus.ko	9 'ka.mjon	7 zlip
CCCV	408	303 'stri.a	68 u.brja.'ko.ne37	'a.trjo	-
CCCVC	295	282 'stril.lo	13 ma'trjos.ka	-	-
CVCC	30	2 post.pran'dja.le-	-	23 'i.noks	5 nord
CCVCC	8	2 'trans.fu.ga	-	4 'du.pleks	2 sport
VCC	2	-	-	-	2 eks
CCCVC	1	-	-	-	1 sprint

• 27.603 words

• 99.260 clusters

• 6.816 unique clusters

the most common:

	intervocalic	initial	final
r	9758 ('ɔ.ro)	1902 (re)	22 (bar)
t	6381 ('da.ta)	873 (te)	11 (bit)
n	6265 ('u.no)	413 (ne)	23 (in)
l	4870 ('bi.le)	674 (lo)	8 (il)
k	2881 ('ɛ.ko)	2359 (ki)	7 (suk)
m	2929 ('i.mo)	1566 (mɔ)	17 (tram)
nt	3623 ('on.ta)	-	2 (sprint)
d	1608 ('ɔ.de)	1490 (di)	1 (po.la'ro.id)
p	1151 ('a.pe)	1603 (per)	4 (dzip)
s	901 ('fu.so)	1492 (su)	27 (bis)
v	1722 ('i.vi)	600 (vi)	-
st	1730 ('us.ta)	330 (sto)	6 (su'dɛst)
b	1164 ('rɔ.ba)	821 (bɔb)	3 (bɔb)
tt	1975 ('ɔt.to)	-	-

the less common:

	intervocalic
bstr	1 (sub'stra.to)
mbrj	1 (em.brjo'na.le)
nglj	1 ('gan.gljo)
nspj	1 (in.spje'ga.bi.le)
rflw	1 (su'per.flwo)
rksj	1 (mar'ksja.no)
rstr	1 (su.per'stra.da)
bbw	1 (ab'bwɔ.no)
ddj	1 (ad'djɛ.tro)
fkj	1 (kaf'kja.no)
grj	1 (ne'grjɛ.ro)
jtj	1 (aj'tja.no)
lkw	1 (al'kwan.to)
...	

• sonority scale

1	complex stops	7	voiced fricatives
2	unvoiced plosives	8	nasals
3	voiced plosives	9	laterals
4	unvoiced affricates	10	trills
5	voiced affricates	11	glides
6	unvoiced fricatives	12	vowels

• 234 unique sonority patterns

examples:

word	son. pattern	contour
sup'erstite	12,10,6,2,12	\\V
rejns'er'ire	12,11,8,6,12	\\V
inkwj'ɛto	12,8,2,11,11,12	\\V-/
ɲ'ɔmo	0,8,12	#//

• 27 unique sonority contours

contour	tokens	example	syll. break
1	V	44254 d'ata, p'atfe	V.CV
2	#/	15812 ke, ʌi	-
3	\\	14485 'onta, v'ispo	VC.CV
4	\\-	8536 k'opto, p'aʌʌa	VC.CV
5	V/	6011 'arja, 'ɔddzi	VC.CV, V.CCV
6	#//	3108 kr'ɛpa, tw'ɔno	-
7	V/	2142 s'empre, m'uskjo	VC.CCV
8	V//	1443 optsj'one, mell'iflwo	VC.CCV, V.CCCV
9	#V	1441 zd'ɛɲɲo, sk'ɛda	-
10	\\-/	937 'ɔppjo, akkl'uzo	VC.CCV
11	#V/	534 str'ega, zgw'ardo	-
12	#-/	232 zv'ago, mnem'ɔniko	-
13	\\#	136 tram, il	-
14	#-//	32 sfr'edzo, zvwotam'ento	-
15	\\#	30 nord, film	-
16	\\V	27 sup'erstite, rejns'er'ire	VC.CCV, VCC.CV
17	\\//	25 ind'ustrja, g'angljo	VC.CCCV
18	#//	17 trjest'ino, zmw'overe	-
19	V-/	16 el'ɔkwjo	V.CCCV
20	V//	11 'ekstra, substr'ato	VC.CCCV
21	\\V/	11 awstr'ale, s'anskrito	VC.CCCV
22	V#	7 faks	-
23	\\-/	4 inkwj'ɛto	VC.CCCV
24	\\-/	3 tr'ansfuga, postkontfij'are	VCC.CV
25	#-/	3 kwj'ɛto	-
26	\\-/	2 akkwjeff'ente	VC.CCCV
27	\\-//	1 postprandj'ale	VCC.CCV