

Structural Lexical Heuristics in the Automatic Analysis of Portuguese

Eckhard Bick

Department of Linguistics, Århus University, Willemoesgade 15 D, DK-8200 Århus N
tel: +45 - 8942 2131, fax: +45 - 86 281397, e-mail: lineb@hum.aau.dk
<http://visl.hum.ou.dk/Linguistics.html>

Abstract

The paper discusses, on the lexical level, the integration of heuristic solutions into a lexicon based and rule governed system for the automatic analysis of unrestricted Portuguese text. In particular, a morphology based analytic approach to lexical heuristics is presented and evaluated. The tagger involved uses a 50.000 entry base form lexicon as well as prefix-, suffix- and inflexion endings lexica to assign part of speech and other morphological tags to every wordform in the text, with recall rates between 99.6% and 99.7%. Multiple readings are subsequently disambiguated by using grammatical rules formulated in the Constraint Grammar formalism. On the next level of analysis, tags for syntactical form and function alternatives are mapped onto the wordforms and disambiguated in a similar way. In spite of using a highly differentiated tag set, the parser yields correctness rates - on running unrestricted and unknown text - of over 99% for morphology/PoS and 97-98% for syntax. A test site with a variety of applications (parsing, corpus searches, interactive grammar teaching and - experimental - MT has been established at <http://visl.hum.ou.dk/Linguistics.html>.

1 Background

In corpus linguistics, most systems of automatic analysis can be classified by measuring them against the bipolarity of rule based versus probabilistic approaches. Thus Karlsson (1995) distinguishes between “pure” rule based or probabilistic systems, hybrid systems and compound systems, i.e. rule based systems supplemented with probabilistic modules, or probabilistic systems with rule based “bias” or postprocessing. As a second parameter, lexicon dependency might be added, since both rules based and probabilistic systems differ internally as to how much use they make of extensive lexica, both in terms of lexical coverage and granularity of lexical information.

The *constraint grammar* (CG) formalism (e.g. Karlsson et al., 1995), which I have been using in my own system¹ for the automatic analysis of unrestricted Portuguese text (Bick, 1996 [1] and 1997 [2]), is both rule governed and lexicon based, focusing on disambiguation of multiply

¹ The system was developed in the framework of a Ph.D.-project at Århus University, over a period of three years, drawing on lexicographic research on Portuguese from an earlier Master's Thesis.

assigned lexical and structural readings as the main tool of analysis. Readings are expressed as sets of word based modular tags. Syntactic structure is covered by using function tags and dependency markers (Bick, 1997 [1]), but I will here concentrate on the lexico-morphological level. Before any Constraint Grammar rules can apply, all (morphologically) *possible* readings have to be identified, and I have to this end developed a preprocessor, that identifies wordforms, polylexical units and sentence boundaries, as well as a morphological analyser for Portuguese using an adapted electronic version of a previously published dictionary (Bick, 1993) in combination with affix- and inflection endings lexica supplemented by corresponding alternation rules for word formation (Bick, 1995). In the analyser's output, every word form is followed by as many tag lines as there are potential readings:

- (1) "<revista>"
 "revista" <+n> <CP> <rr> N F S
 "revestir" <vt> <de^vtp> <de^vrp> V PR 1/3S SUBJ VFIN
 "revistar" <vt> V IMP 2S VFIN
 "revistar" <vt> V PR 3S IND VFIN
 "rever" <vt> <vi> V PC

With a CG-term, such an ambiguous list of readings is called a *cohort*. In the example, the word form 'revista' has one noun-reading (female singular) and four (!) verb-readings, the latter covering three different base forms, subjunctive, imperative, indicative present tense and participle readings. Conventionally, PoS and morphological features are regarded as primary tags and coded by capital letters. In addition there can be secondary lexical information about valency and semantical class, marked by <> bracketing.

A *constraint grammar rule* brings the ambiguity problem to the foreground by specifying which reading (out of a cohort of ambiguous readings for a given word) is impossible (and thus to be discarded) or mandatory (and thus to be chosen) in a given sentence-context. For instance, a rule might discard a finite verb reading after a preposition (2a) , or when another - unambiguous - finite verb is already found in the same clause, with no coordinators present (2b).²

- (2a) REMOVE (VFIN) IF (-1 PRP)
 [discard finite verb readings (VFIN) if the first word to the left (-1) is a preposition PRP]
- (2b) REMOVE (VFIN) IF (*1C VFIN BARRIER CLB OR KC) (NOT *-1 CLB-WORD)
 [discard VFIN if there is another unambiguous (C) finite verb (VFIN) anywhere to the right (*1) with no clause-boundary (CLB) and coordinating conjunction (KC) interfering (BARRIER). Discard only if there is no subordinator (CLB-WORD) anywhere to the left (*-1)]

² Ordinarily, this disambiguation process works on whole cohort lines, i.e. distinguishes between PoS, base form and inflection, but tolerates competing valency options. However, on a higher level of analysis, I have introduced valency and semantical disambiguation, too. This can be very useful for polysemy resolution, like in "rever", where the transitive <vt> - intransitive <vi> distinction has a meaning correlate: 'tornar a ver' [see again] vs. 'transudar' [leak through]. Likewise, "revista" followed by a name <+n> or being read (semantical class <rr>) is more likely to be a newspaper than an inspection (semantical class <CP> for action: +CONTROL, +PERFECTIVE).

With current software, before an analysis run, all rules are translated into a finite state network by a compiler program, yielding the actual parser. The Portuguese grammar was originally written in the formalism suggested by Pasi Tapanainen's first compiler implementation, but later rewritten to match the notation used in his new CG-2 parser compiler (Tapanainen, 1996).

By applying the rule set several times, the parser renders more and more words in the sentence unambiguous, and in the end, only one reading is left for every word. Since the individual rule can be made very "cautious" by adding more context conditions, and since the last surviving reading will never be discarded, the formalism is very robust. Even imperfect input will yield *some* parse. Unlike probabilistic systems, where "manual interference" as in the introduction of bias on behalf of irregular phenomena often has an adverse side-effect on the overall performance of the parser (due to interference with the ordinary statistical "rules" based on the *regular* "majority" phenomena), Constraint Grammar tolerates and even encourages the incremental "piecemeal" addition of exceptions and context conditions for individual rules (For a comparison of statistical and constraint-based methods see Chanod & Tapanainen, 1994).

2. System Performance

If they can be made to work on free text, rule based systems can achieve very low error rates. While state-of-the-art probabilistic taggers still have error rates of over three percent³, even for PoS tagging, CG based systems fare somewhat better. For English word class error rates of under 0.3% have been reported at a disambiguation level of 94-97% (Voutilainen, 1992). For my own Portuguese CG system, test runs runs with near 100% disambiguation on fiction and news texts suggest a correctness rate of over 99% for morphology and part of speech, when analysing unknown unrestricted text⁴. For syntax the figures are 98% for classical literary prose (Eça de Queiroz, "O tesouro") and 97% for the more inventive "journalese" of news magazine texts (VEJA,9.12.1992),as shown in table(3):

³ Compare, for English, (Garside et.al., 1987) on the HMM based CLAWS system, (Francis and Kucera, 1992) on recovering PoS tags from the Brown corpus, Ratnaparkhi's maximum-entropy tagger trained on the Penn Treebank (Marcus et al., 1993) or Brill's stochastic tagger using automated learning (Brill, 1992). For German the Morphy system described in (Lezius et. al., 1996) achieved an accuracy of 95.9%.

⁴ The test texts used were not part of the benchmark corpus used to develop the rules, and fresh text chunks were used for every new test. The present grammar, however, still being improved, does incorporate changes made as a result of test run errors.

(3) System performance on the PoS and syntactic levels:

Text: Error types:	<i>O tesouro</i> ca. 2500 words		<i>VEJA 1</i> ca. 4800 words		<i>VEJA 2</i> ca. 3140 words	
	errors	correct- ness	errors	correct- ness	errors	correct- ness
Part-of-speech errors	16		15		24	
Base-form & flexion errors	1		2		2	
All morphological errors	17	99.3 %	17	99.7 %	26	99.2 %
syntactic: word & phrases	54		118		101	
syntactic: subclauses	10		11		13	
All syntactic errors	64	97.4 %	129	97.3 %	114	96.4 %
"local" syntactic errors due to PoS/morphological errors	- 27		- 23		- 28	
Purely syntactic errors	37	98.5 %	106	97.8 %	86	97.3 %

3. Lexico-morphological heuristics

Yet even in a rule based CG system, heuristics can be quite useful (for English, see Karlsson et. al., 1995). Thus rules are usually grouped according to their "safety", i.e. their statistical tendency to make errors. Less safe rules can be added as a heuristic level on top of a kernel of safe rules, and will be applied *after* these. Also, statistical inspired "rarity tags" (<Rare>) can be added to certain less probable readings in the lexicon, and then referred to by contextual disambiguation rules. A third field for the application of heuristics is on the *analyser level*, i.e. concerns the (lexico-morphological) *input* of the disambiguation rule system. It is this third type of heuristics I am concerned with here.

Since the higher levels of the parsing system (for example, PoS and syntax) are technically rule based disambiguators, they need *some* reading for every word to work on, which is why even unanalyzable word forms (i.e. word forms that can not be reduced to a root found in the analyser's lexicon) need to be given one or more heuristic readings with regard to word class and flexion morphology. The majority of such cases is accounted for by unknown proper nouns (1-2% of all words, depending on text type), and can be handled by assigning heuristic PROP tags to *all* capitalised words in certain contexts, and then *adding* any competing *analytical* analysis, especially in the case of sentence initial position, leaving the final decision to the rule based disambiguation module. This way, though 80% of all nouns in my corpus need heuristical treatment, the error rate for the PROP class as a whole can be kept at 2%, not too far from the taggers overall PoS error rate (< 1%).

3.1 "Unanalyzable⁵ words": typology and statistics

Though accounting for only 0.3-0.5% of word forms in running text, other types of analysis failures (i.e. word forms that can not be reduced to a root found in the analyser's lexicon) are more difficult to handle, due to their functional diversity and the lack of a clear morphological marker. Table (4) provides an error typology for a 131.981 word literature and secondary literature corpus (The RNP depository of Brazilian literature), containing 604 unanalyzable words in the test run.

Three main groups may be distinguished, comprising of roughly one third of the cases each:

- a) orthographical errors (shaded in the table, and partially corrected before heuristics proper by an accent module recognising regular regional spelling variations)
- b) unknown and underivable Portuguese words or abbreviations
- c) unknown foreign loan words

(4) Error types in "unanalyzable" words:

DOMAIN	NUMBER OF TOKENS	PERCENTAGE	
Foreign	232	38.4	
orthographic variation (European/accenuation)	125	20.7	<i>Correctables</i>
other port. Orthographic	74	12.3	<i>Misspellings</i>
non-capitalised names and abbreviations	37	6.1	<i>Encyclopaedic lexicon failures</i>
names and name roots	18	3.0	
abbreviations	19	3.1	
root not found in lexicon	119	19.7	<i>Core lexicon failures</i>
found in Aurelio ⁶	91	15.1	
not found in Aurelio	28	4.6	
derivation/flexion problem	15	2.5	<i>Affix lexicon failures</i>
suffix	8	1.3	
prefix	3	0.5	
flexion ending	2	0.3	
alternation information	2	0.3	
other	2	0.3	
SUM	604	100.0	

3.2 Analytical morphological heuristics

For optimal performance, the three groups mentioned above would require different strategies. Foreign words appearing in running Portuguese text are typically nouns or noun phrases, and

⁵ In this paper I intend "unanalysable word forms" to mean word forms that cannot - by derivation and/or inflexional analysis - be reduced to a root found in the analyser's lexicon. Of course, only part of these - typing errors and foreign language quotes - are *really* unanalysable, while others might be covered by enlarging the lexicon or enhancing the scientific derivation list.

⁶ "Novo Dicionário Aurelio" is the largest monolingual dictionary of Brazilian Portuguese.

trying to identify verbal elements only causes trouble. In "real" Portuguese words without spelling errors, structural clues - like flexion endings and suffixes - should be emphasised. These will be meaningful in misspelled Portuguese words, too, but, in addition, specific rules about letter manipulation (doubling of letters, missing letters, letter inversion, missing blanks etc.) and even knowledge about keyboard characteristics might make a difference.

Motivated by a grammatical perspective rather than probabilistics, my approach has been to emphasise groups (a) and (b) and look for *Portuguese* morphological clues in words with unknown stems. Since prefixes have very little bearing on the probability of a word's word class or flexional categories, only the flexion endings and suffix lexica are used. As it also does in ordinary morphological analysis, the tagger tries to identify a word from the right, i.e. backwards, cutting off potential endings or suffixes and checking for the remaining stem in the root lexicon (the main lexicon). Normally, for (multiply) *analysed* words, using Karlsson's law (Karlsson 1992, 1995)⁷, the Portuguese analyser would try to make the root as long as possible, and to use as few derivational layers⁸ as possible. For (system-internally) *unanalyzable* words, however, I use the opposite strategy: Since I am looking for a hypothetical root, flexion endings and suffixes are all I've got, and I try to make their half of the word (the right hand part) as large as possible.

Working with a minimal root length of 3 letters, and calling my hypothetical root 'xxx', I will start by replacing only the first 3 letters of the word in question by 'xxx' and try for an analysis, then I will replace the first 4 letters by 'xxx', and so on, until - if necessary - the whole word is replaced by 'xxx'.⁹ For a word like *ontogeneticamente* the rewriting record will yield the chain below. Here, the full chain is given, with all readings it would encounter on its way. In the real case, however, the tagger - preferring long derivations/endings to short ones - would stop searching at the *xxxticamente* -level, where the first group of readings is found. In fact, the adverbial use of an adjectively suffixed word is much more likely than hitting upon, say, a "root-only" noun whose last 9 letters happen to include both the '-ico' and the '-mente' letter chains by chance.

(5) *ontogeneticamente* -> no analysis

xxxogeneticamente
xxxgeneticamente
xxxeneticamente
xxxneticamente

⁷ Karlsson's law states, that of two morphological analyses of different derivational complexity, the one with *fewer* elements is almost always the correct one.

⁸ Karlsson's law can be applied to any string of free (i.e. compounding), derivational or inflexional morphemes, but the frequency of ambiguity types with respect to these three elements will differ from language to language - thus, in Portuguese, compounding is much rarer than in most Germanic languages, while Swedish, the language for which Karlsson's law was originally formulated, does have compounding, but not as rich an inflexion morphology.

⁹ A similar method of partial morphological recognition and circumstantial categorization might be responsible for a human being's successful inflectional and syntactic treatment of unknown words in a known language; the Portuguese word games "collorido" (president Collor & *colorido* - 'coloured') and "tucanagem" (the party of the *tucanos* & *sacanagem* - 'dirty work'), for instance, will not be understood by a cultural novice in Brazil, even if he is a native speaker of European Portuguese - but he will still be able to identify both as singular, the first as a past participle ('-do') and the second as an abstract noun ('-agem') of the feminine gender.

<i>xxxeticamente</i>	
<i>xxxiticamente</i>	-> suffix '-ico' (variation '-tico') + adverbial ending '-mente' "ontogene" <DERS -ico [ATTR]> <deadj> ADV
<i>xxxicamente</i>	-> suffix '-ico' + adverbial ending '-mente' "ontogene" <DERS -ico [ATTR]> <deadj> ADV
<i>xxxicamente</i>	
<i>xxxamente</i>	-> adverbial ending '-mente' (variation '-amente') "ontogenetico" <xxxo> <deadj> ADV
<i>xxxmente</i>	
<i>xxxente</i>	-> "present participle"-suffix '-ente' "ontogeneticamer" <DERS -ente [PART.PR]> ADJ M/F S "ontogeneticamer" <DERS -ente [AGENT]> N M/F S -> causative suffix '-entar' ¹⁰ + verbal flexion ending '-e' "ontogeneticam" <DERS -ar [CAUSE]> V PR 1/3S SUBJ VFIN
<i>xxxnte</i>	
<i>xxxte</i>	
<i>xxx</i>	-> verbal flexion ending '-e' "ontogeneticamenter" <xxxer> V IMP 2S VFIN "ontogeneticamentir" <xxxir> V IMP 2S VFIN ### "ontogeneticamenter" <xxxer> V PR 3S IND VFIN "ontogeneticamentir" <xxxir> V PR 3S IND VFIN ### "ontogeneticamentar" <xxxar> V PR 1/3S SUBJ VFIN
<i>xxx</i>	-> no derivation or flexion "ontogeneticamente" <xxx> N F S "ontogeneticamente" <xxx> N M S

Roots with 'xxx' are present in the core lexicon alongside the "real" roots, including the necessary stem alternations¹¹ for verbs (here, BbCc for different root-stressed forms and AaiD for endings-stressed forms):

(6)

root	word class	alternation subclass	lexeme ID	target of analysis
xxx	<sm>		54573	masculine noun, typically foreign
xxxar	<amf>		59547	Portuguese '-ar'-adjective*
xxxar-	<vt>	AaiD	54578	endings-stressed forms of '-ar'-verbs
xxxer	<sm>		54666	masculine noun, typically English*
xxxia	<sf>		54665	feminine noun, Latin-Portuguese*
xxxo	<sm>		54582	masculine noun, typically Portuguese

¹⁰ This suffix is regarded as a variant of '-ar', and therefore normalized in the DER-tag: <DERS -ar [CAUSE]>.

¹¹ Here, BbCc for different root-stressed forms and AaiD for endings-stressed forms, with D, for example, meaning a root to be used with future subjunctive endings.

Besides the typical stems ending in '-o', '-a' and '-r', default stems consisting of a plain 'xxx' have been entered to accommodate for foreign nouns with "un-Portuguese" spelling. Like many other languages, Portuguese will force its own gender system even onto foreign loan words, so a masculine and a feminine case must be distinguished, for later use in the tagger's disambiguation module.

Since the analyser's heuristics for unknown words prefers readings with endings (or suffixes) to those without, and longer ones to shorter ones, verbal readings (especially those with inflexion morphemes in 'r', 'a' or 'o') have a "natural" advantage over what really should be nouns or adjectives, especially when these appear in their uninflected singular base form. Lexicon-wise, this tendency is countered by adding three of the most commonly ignored nominal cases specifically into the lexicon: (a) English '-er' nouns otherwise only taken as Portuguese infinitives, (b) Latin-Portuguese '-ia' nouns otherwise only read as verbal forms in the imperfeito tense, and (c) '-ar' adjectives otherwise analysed only as infinitives.

Rule-wise, verbal readings alone are not allowed to stop the heuristics-machine, it will proceed until it finds a reading with another word class. So, the process is set to ignore *verbal* readings on its way down the chain of hypothetical word forms with ever shorter suffix/endings-parts. Thus, the heuristics-machine will *record* verbal readings, but only stop if a noun, adjective or adverb reading is found in that level's cohort (list of readings). In this context, participles and gerunds - though verbal - are treated as "adjectives" and "adverbs", respectively, because they feature very characteristic endings ('-ado', '-ido', '-ando', '-endo', '-indo').

This raises the possibility of the heuristics-machine progressing from multi-derived analyses (with one or more suffixes) to simple analyses (without suffixes) before it encounters a non-verbal reading. In this case, the application of Karlsson's law does still make sense, and when the heuristics-machine hands its results over to the local disambiguation module, this will select the readings of lowest derivational complexity, weeding out all (read: verbal!) readings containing more (read: verbal!) suffixes than the group selected. In the misspelled French word '*entaente*', for example, the verbal reading:

(7a) "enta" <DERS -(ent)ar [CAUSE]> V PR 1/3S SUBJ VFIN,
from the 'xxxente'-level, is removed, leaving only underived verbal readings - from the 'xxx'-level, along with the desired noun singular reading from the 'xxx'-level.

(7b) entaente ALT xxxaente ALT xxxe ALT xxx
 "entaenter" <xxxer> V IMP 2S VFIN
 "entaenter" <xxxer> V PR 3S IND VFIN
 "entaentar" <xxxar> V PR 1/3S SUBJ VFIN
 "entaente" <xxx> N F S
 "entaente" <xxx> N M S

Since all disambiguation not related to Karlsson's law is referred to the CG-module, the word class choice between V and N will be contextual (and rule based), as well as the morphological sub-choice of mode (IMP - PR) for the verb, and gender (M - F) for the noun. In the prototypical case of a preceding article, the verb reading is ruled out by:

(8a) REMOVE (V) IF (-1 ART)
and the gender choice is then taken by agreement rules such as:

- (8b) REMOVE (N M) IF (- 1C DET) (NOT -1 M)
 REMOVE (N F) IF (- 1C DET) (NOT -1 F)

Consider the following examples of "unanalyzable" words from real corpus sentences, where the final output, after morphological contextual disambiguation, is given:

- (9a) inventimanhas ALT xxxas (also: one ADJ and three rare V-readings)
 "inventimanha" <xxx> N F P 'tricks'
- (9b) araraquarenses ALT xxxenses (3 other ADJ readings removed by local disambiguation)
 "araraquar" <DERS -ense [PATR]> <jh> <jn> ADJ M/F P 'from Araraquara'
- (9c) sombrancelhas ALT xxxas (also: one ADJ and three rare V-readings)
 "sombrancelha" <xxx> N F P '=sobrancelhas - eye brows'
- (9d) cast ALT xxx (also: N F S)
 "cast" <*1> <*2> <xxx> N M S 'English: cast'

In (9a) and (9b) the parser assigns correct readings to unknown, but wellformed Portuguese words. Depending on the orthodoxy of the fusion process, these affixes may be recognised (9b), or not (9a). That affix recognition is important, can be seen from the fact that all competing analyses in (9b) - but not in (9a) - have the correct PoS tag. What is special about (9c), is the (phonetical?) misspelling ('sombrancelhas') of an otherwise ordinary Portuguese word. Even so, with the help of the surviving morphological clues and contextual disambiguation, the parser is able to assign the right analysis in most cases, especially if the words still look Portuguese. (9d), finally, is the hard case - foreign loan words. English 'cast' does not fit with any Portuguese flexion ending, therefore the default reading N is assigned, gender disambiguation relying on NP-context.

In order to test the parser's performance and to identify the strengths and weaknesses of the heuristics strategy of the parser, I have manually inspected 757 "running" instances¹² of lower case word forms where the parser's disambiguation module received its input from the tagger's heuristics module. The first column shows the word class analysis chosen, and inside the three groups (errors, Portuguese, foreign) the left column gives the number of correct analyses, whereas the right column offers statistics about the mistakes, specifying - and quantifying - what the analysis *should* have been.

¹² The words comprise all "unanalysable" word forms in my corpus, that begin with the letters 'a' and 'b'. Since the relative distribution of foreign loan words and Portuguese words depends on which initial letters one works on ('a', for one, is over-representative of Portuguese words, whereas 'x', 'w' and 'y' are English-only domains), no conclusions can be drawn about these two groups' relative percentages. Inside the Portuguese group, however, the distribution between real words and misspellings may be assumed to be fairly alphabet-independent. Any way, the sampling technique has no significance for error frequencies or distribution in relation to word class, which was the main objective in this case.

(10) Word class distribution and parser performance in "unanalyzable" words (VEJA news text)

analysis	A) orthographical errors		B) Portuguese words		C) foreign words ¹³		all	
	correct	other	correct	other	correct	other	correct	other
N	119	ADJ 8 ADV 8 VFIN 3 PRON 1 DET 1 PRP 1	212	ADJ 3	226	ADV 11 ADJ 3 PRON 2 PRP 2	557	43
ADJ	25	N 8 GER 2	95	N 7	8	-	128	17
ADV	3	-	5	-	-	-	8	-
VFIN	13	N 4 PCP 1 ADV 1	9	N 4 ADJ 2	-	N 7 ADJ 1	22	20
PCP	10	-	16	-	-	-	26	-
GER	3	-	-	-	-	-	3	-
INF	9	-	4	-	-	N 4	13	4
	182	38	341	16	234	30	757	84
		(17.3%)		(4.5%)		(11.4%)		(10.0%)

The table shows that, when using lexical heuristics, the parser performs best - not entirely surprisingly - for wellformed Portuguese words (B). Of 323 nouns and adjectives in group B, only 16 (5%) were misanalysed as false positives or false negatives. The probability for an assigned N-tag being correct is as high as 98.6%, for the underrepresented adverb and non-finite verbal class even 100%. All false positive nominal readings (N and ADJ) are still in the nominal class, a fact that is quite favourable for later syntactic analysis.

Figures are lower for group C, unknown loan words, where the chance of an N-tag being correct is only 92.6%, even when allowing for a name-chain-like N-analysis of English adjectives integrated in noun clusters of the type 'big boss'. Finite verb readings, though rare (due to lacking flexion indicators), are of course all failures, and only the little adjective group was a hit, the few cases being triggered by morphologically "Portuguesish" Spanish or Italian words.

The results in group A (misspellings) resemble those of group B, with a good performance for classes with clear endings, i.e. non-finite verbs and '-mente'-adverbs, and a bad performance for finite verb forms. For the large nominal groups figures are somewhat lower: 84.4% of N-tags, and only 71.4% of ADJ-tags are correct - though most false positive ADJ-tags are still within the nominal range. The lower figures can be partly explained by the fact that misspelled closed class words (adverbs, pronouns and the like) will get the (default, but wrong) noun reading - a technique that works somewhat better and more naturally for foreign loan words (C), which often are "terms" imported together with the thing or concept they stand for, or names. Also, the

¹³ Only individual words and short integrated groups are treated, foreign language sentences or syntactically complex quotations are treated as "corpus fall-out" in this table.

percentage of "simplex"¹⁴ words without affixes is much higher among the misspellings in group A than in group B, where all simplex words - being spelled correctly - would have been recognised in the lexicon anyway, due to the good lexicon coverage *before* getting to the heuristics module. Therefore, nouns and adjectives in group A lack the structural information of suffixes that helps the parser in group B: 'xxxo' looks definitely less adjectival than 'xxxístico'. In particular, 'xxxo' invites the N/ADJ-confusion, whereas many suffixes are clearly N or ADJ. Thus, '-ístico' yields a safe adjective reading.

4. Special - "deviant" - word class probabilities for the heuristics module

Is it possible, apart from morphological-structural clues, to use "probabilistics pure" for deciding on word class tags for "unanalyzable" words? In order to answer this question, I will - in table (14) - rearrange information from table (13) and compare it to whole text data (in this case, from a 197.029 word stretch of the mixed genre Borba-Ramsey corpus). Here, I will only be concerned with the open word classes, nominal, verbal and '-mente'-adverbial.

(11) Open word class frequency for "unanalyzable" words as compared to whole text figures

	whole text		"unanalyzable" words						
			orthographical errors		Portuguese words		foreign words		all heuristics
analyse s	%	cases	%	cases	%	cases	%	cases	%
N	47.38	131	63.59	232	63.39	237	95.18	600	73.08
ADJ	12.79	33	16.02	100	27.32	12	4.82	145	17.66
ADV ¹⁵	1.26	3 (+9)	1.46	5	1.37	- (+11)	-	8	0.97
VFIN	24.96	16	7.77	9	2.46	-	-	25	3.05
PCP	4.96	11	5.34	16	4.37	-	-	27	3.29
GER	2.47	3	1.46	-	-	-	-	3	0.37
INF	6.17	9	4.37	4	1.09	-	-	13	1.58
		206		366		249		821	

Among other things, the table shows that the noun bias in "unanalyzable" words is much stronger than in Portuguese text as a whole, the difference being most marked in foreign loan words. The opposite is true of finite verbs which show a strong tendency to be analysable. Finite verbs are virtually absent from the unknown loan word group. For the non-finite verbal classes the distribution pattern is fairly uniform, again with the exception of foreign loan words.

¹⁴ "Simplex" words are here defined as words that can be found in the root lexicon without prior removal of prefixes or suffixes. Of course, the larger the lexicon the higher the likelihood of an (etymologically) affix-bearing word appearing in the lexicon, - and thus not needing "live" derivation from the parser.

¹⁵ Only deadjectival '-mente'-adverbs can meaningfully be guessed at heuristically, and therefore only they should enter into the statistics for word class guessing. Also the base line figure of 1.26% for normal text is for '-mente'-adverbs only, the overall ADV frequency is nearly 12 times as high. Since non-'mente'-adverbs are a closed class in Portuguese, the latter will be absent from the heuristics class of wellformed unknown Portuguese words, but in the foreign loan word group and the orthographical error group they will appear in the false positive section of other word classes (numbers given here in parentheses). In the orthographical error group, both '-mente'-adverbs and closed class adverbs can occur, the first as correct ADV-hits, the other usually as false positive nouns (for instance, 'aimda').

As might be expected, among the "unanalyzable" words, orthographical errors and correct Portuguese words show a remarkably similar word class distribution.

A lesson from the above findings might be to opt for noun readings and against finite verb readings in "unanalysable" words, when in doubt, especially where no Portuguese flexion ending or suffix can be found, suggesting foreign material. As a matter of fact, this strategy has since been implemented in the system, in the form of heuristical disambiguation rules, that discard VFIN readings and chose N readings for <MORF-HEUR> words, where lower level (i.e. safe) CG-rules haven't been able to decide the case contextually.

5. Conclusion

It can be shown that lexico-morphological heuristics - at least for a morphology-rich language like Portuguese - can be based on structural clues and the systematic exploitation of derivational and inflectional sublexica. Applied to improve analyser recall on the input level of a Constraint Grammar system, the described technique positively contributed to the overall performance of a lexicon based rule governed tagger/parser. Correctness rates of more than 99% were achieved for the morphological/PoS tagger module, with heuristic error rates running at 2% for proper name heuristics and 4.5% for the heuristical analysis of other unrecognised, but correctly spelled Portuguese word forms. In all, heuristic analysis was needed for 80% of all proper nouns (amounting to ca. 2% of running word forms in news text), but for less than 0.4% of non-name word forms. Finally, word class frequency counts suggest that PoS probabilities for "unanalyzable" words in Portuguese texts are quite different from those for the language on the whole.

References

Bick, Eckhard, *Portugisisk - Dansk Ordbog*, Mnemo, Århus, 1993, 1995, 1997.

Bick, Eckhard, *The Parsing System "Palavras", Documentation*, unpublished Ph.D. project evaluation, 1995, 1997.

Bick, Eckhard, "Automatic Parsing of Portuguese", in *Proceedings of the Second Workshop on Computational Processing of Written Portuguese*, Curitiba, 1996.

Bick, Eckhard, "Dependensstrukturer i Constraint Grammar Syntaks for Portugisisk", in: Brøndsted, Tom & Lytje, Inger (eds), *Sprog og Multimedier*, Aalborg, 1997.

Bick, Eckhard, "Automatisk analyse af portugisisk skriftsprog", in: Jensen, Per Anker & Jørgensen, Stig. W. & Hørning, Anette (eds), *Danske ph.d.-projekter i datalingvistik, formel lingvistik og sprogteknologi*, pp. 22-20, Kolding, 1997.

Brill, Eric, "A Simple Rule-based Part of Speech Tagger", in *Proceedings of the Third Conference on Applied Natural Language Processing*, ACL, Trento, Italy, 1992.

Chanod, Jean-Pierre & Tapanainen, Pasi, "Tagging French - comparing a statistical and a constraint-based method", adapted from: *Statistical and Constraint-based Taggers for French*, Technical report MLTT-016, Rank Xerox Research Centre, Grenoble, 1994.

Francis, W.N. & Kucera, F., *Frequency Analysis of English Usage*, Houghton Mifflin, 1982.

Garside, Roger & Leech, Geoffrey & Sampson, Geoffrey (eds.), *The Computational Analysis of English. A Corpus-Based Approach*, London, 1987.

Karlsson, Fred, "SWETWOL: A Comprehensive Morphological Analyser for Swedish", in *Nordic Journal of Linguistics* 15, 1992, pp. 1-45.

Karlsson, Fred & Voutilainen, Atro & Heikkilä, Juka & Anttila, Arto (eds.), *Constraint Grammar, A Language-Independent System for Parsing Unrestricted Text*, Mouton de Gruyter, Berlin 1995.

Lezius, Wolfgang & Rapp, Reinhard & Wettler, Manfred, "A Morphology-System and Part-of-Speech Tagger for German", in: Dafydd Gibbon (ed.): *Natural Language Processing and Speech Technology*, Berlin, 1996.

Marcus, Mitchell, "New trends in natural language processing: Statistical natural language processing", paper presented at the colloquium *Human-Machine Communication by Voice*, organized by Lawrence R. Rabiner, held by the National Academy of Sciences at *The Arnold and Mabel Beckman Center* in Irvine, USA, Feb. 8-9, 1993.

Tapanainen, Pasi, "The Constraint Grammar Parser CG-2", University of Helsinki, Department of Linguistics, Publications no. 27, 1996.

Voutilainen, Atro & Heikkilä, Juka & Anttila, Arto, *Constraint Grammar of English, A Performance-Oriented Introduction*, Publication No. 21, Department of General Linguistics, University of Helsinki, 1992.