

# Supervised Word Sense Disambiguation for Venetan: a Proof-of-Concept Experiment

Costanza Conforti and Alexander Fraser

Center for Information and Language Processing  
Ludwig-Maximilians University of Munich  
Oettingenstrae 67, 80583 - Munich, Germany

## Abstract

Word Sense Disambiguation (WSD) is a classification task that consists of determining which of the senses of an ambiguous word is activated in a specific context. Research in this field has primarily concentrated on investigating English and a few other well-resourced languages. Recently, studies done on a corpus of Old English (Wunderlich 2015) showed that, even with limited resources, it is still possible to approach the problem of WSD. In this paper a WSD system has been developed for the Low Resource Language (LRL) Venetan, which has recently received some attention from the Natural Language Processing (NLP) community. Our main contributions are twofold: first, we select and annotate a corpus for Venetan, considering two words (one abstract and one concrete term) and using two levels of annotation (fine- and coarse-grained), reporting on annotator agreement. Second, we report results of proof-of-concept experiments of supervised WSD performed with Support Vector Machines on this corpus. To our knowledge, our work is the first time that WSD for a European Dialect like Venetan has been studied.

## Introduction

NLP research in the field of LRLs is hindered by a number of factors, like the lack of lexicographical resources, combined with the scarcity of available and labeled data. NLP for Italian dialects<sup>1</sup> has not been investigated much. Even if they are usually relatively close to the standard language, adapting existing tool for Italian can be very challenging, as dialects are very fragmented and show many orthographic and morphological variants. This constitutes a general challenge in NLP applications focusing on LRL and Old Languages. Aside from a few contributions (Bortolotti 2005), work on Italian dialects mainly concentrated on Venetan, a LRL primarily spoken in the Northeastern regions of Italy. The term LRL can be used to refer to a number of different situations. Following the classification proposed by Singh

Copyright © 2017, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

<sup>1</sup>In this paper, we use the word *language* and *dialect* indiscriminately. Partly autonomous, Italian dialects came into being more or less in the same period through transformations of Latin. Many centuries later, one of them, the dialect of Florence, became the official language of the Italian State in 1861. However, from a historical and linguistic point of view, each dialect can be seen as a language on its own (Berruto 2005).

(2008), Venetan can be considered a resource scarce language, as it has been widely studied from a linguistic point of view, but very few resources and tools for NLP are available. In recent years, the research community started to investigate some classic NLP tasks like morphological analysis (Tonelli 2010) and POS tagging, (Jaber 2011), including a preliminary study on Statistical Machine Translation from Venetan into English (Delmonte 2009).

## A brief description of Venetan

Venetan is a Gallo-Italic language spoken in the Northeast region of Italy, where it is spoken as mother tongue by about 3.852.500 people. Ethnologue reports also 50000 speakers in Croatia and about 4 million in Brazil, where this language is called *Talian*<sup>2</sup>. Venetan is in vigorous use: proficiency among local speakers reaches 75% of the population and it is widely spoken within generations (Tonelli 2010). A number of newspapers are published partially in Venetan and some radio stations broadcast in Venetan; moreover, a small community of Venetans is very active on the Internet, where many blogs and websites can be found, including a version of Wikipedia<sup>3</sup>. Far from constituting a standardized language, Venetan constitutes a regional continuum. Linguists have identified at least four main varieties of Venetan, which show peculiarities in morphology, syntax and lexicon but are still mutually comprehensible. Therefore, Venetan can be considered as a diasystem, where speakers use their own variety in everyday life and manage to understand each other (Tonelli 2010). Venetan is usually written with the Italian alphabet, plus some special characters. Some attempts to unify orthography have been made, like the project promoted by the Veneto Region that culminated in the *Manual of Venetian Orthography*. However, these guidelines are far from being universally accepted and inconsistencies in the orthography and strong morphological variations are both frequent.

## Corpus creation

In this Section, we describe the linguistic resources and the preprocessing pipeline used to extract candidates for WSD and for annotating the training samples.

<sup>2</sup><https://www.ethnologue.com/language/vec>

<sup>3</sup><https://vec.wikipedia.org/wiki/V%C3%A8neto>

## Resources Description

We used the STILVEN corpus as our source for extracting training examples. STILVEN is a project founded by the Veneto Region in 2008<sup>4</sup> and carried out by researchers from the Ca' Foscari University (Venice) and IRST-FBK (Trento). The project involved the implementation of a morphological analyzer for Venetan which permitted the creation of a corpus with a homogenized orthography (Tonelli 2010), and the development of a Venetan POS-tagger (Jaber 2011). To our knowledge, STILVEN is the only available corpus for an Italian dialect which is orthographically normalized and POS-tagged. The corpus collects very heterogeneous texts, including children stories, famous quotes, a manual of Venetan orthography rules and translations of a book about American history and of *The Little Prince* by Antoine de Saint-Exupry. Statistical information about the STILVEN corpus are listed in Table 1.

Token count	133734
Type count	13321
Ratio of (token count/type count)	10.03
Total number of sentences	13058
Average sentence length	10.24
Minimum sentence length	2
Maximum sentence length	77

Table 1: STILVEN Corpus statistics

## Candidate Extraction

In order to obtain a list of polysemous words, we searched in the STILVEN corpus for common nouns with minimum token frequency of 95 and minimum word length of 3 characters<sup>5</sup>. Seven words meet these criteria. All seven terms were, at different levels, polysemous<sup>6</sup> (Table 2). In the following sections, the annotation process and WSD results for the words *parte* (Engl. *part, somewhere*) and *omeni* (Engl. *men, soldiers*) are analysed in detail. Concerning the occurrences of these two words, we excluded the sentences where the candidate was part of a collocation (as for example *omeni de afari*, Engl. *businessmen*). Three occurrences of *parte* were also discarded, as they had been wrongly tagged as noun while they were actually a verb (Eng. *leave*). The final occurrence count was of 98 occurrences for the word *omeni* and 121 for *parte*.

<sup>4</sup>[project.cgm.unive.it/stilven\\_en.html](http://project.cgm.unive.it/stilven_en.html)

<sup>5</sup>We considered only words occurring with a frequency higher than 95 in order to obtain a corpus sufficiently large to be used for training a WSD system. This methodology is similar to the one applied in (Wunderlich 2015).

<sup>6</sup>However, three words in this list have been discarded due to different reasons. The sense distributions of *idea* (Engl. *idea*) and *man* (Engl. *help, hand*) were too unbalanced (90/11/106 considering the three senses of *idea*, 113/23 considering the two senses of *man*). The word *dito*, (Engl. *finger, proverb*) had been tagged as a noun, but in more than a half of the occurrences it was actually the past participle of the verb *say*.

Word Token	Count	English Translation
dito	251	<i>finger, proverb</i>
idea	208	<i>idea</i>
roba	208	<i>stuff, thing, food</i>
man	138	<i>help, hand</i>
parte	144	<i>part, somewhere</i>
tenpo	113	<i>time</i>
omeni	99	<i>men, soldiers</i>

Table 2: Candidates for WSD from the STILVEN corpus

## Synset Definition and Candidate Annotation

To generate training data for supervised classification, all occurrences of *parte* and *omeni* were manually labeled with the sense activated in the sentence. To our knowledge, no lexicographic resource is available for Venetan. Therefore, we proceeded as follows: first of all, we looked up in the Venetan-Italian translation dictionary *El Galepin*<sup>7</sup> for the Italian translations of each Venetan word. Then, the Italian and English WordNet<sup>8</sup> were consulted in order to collect the *synsets* of each translation. With the help of a Venetan native speaker, we merged the most related and fine-grained *synsets* in order to obtain a small final set of clearly distinct senses. We selected three senses for the word *omeni*, while for the word *parte*, we consider three coarse-grained senses and six fine-grained senses (see Table 3). Finally, the set of senses of each word, with the related corpus, has been given to two non-professional annotators separately<sup>9</sup>.

The effort of manual annotation was considerable, as sense annotation is a very difficult undertaking (Wilks 1998). As the annotation sessions with the second judge were quite laborious, we redefined the task from a classification to a discrimination task, as was done by Gale (1992): the second annotator received a sentence labeled by the first judge and had to report whether she agreed or not with the classification. For the final annotation, when the annotators disagreed, the label proposed by the first judge was taken as the gold standard. As reported in Table 3, the sense distributions of *omeni* and *parte* with coarse-grained labels are quite unbalanced, whereas the fine-grained classification of *parte* is more uniform. Judges' agreement for *parte* using fine-grained labels was lower than using coarse-grained labels (~79.0% vs. ~93.9%). Judges' agreement for the concrete word *omeni* was lower than for the abstract term *parte* (~89.7%). These values are far from the 96.8% agreement obtained by Gale (1992), but this difference can be explained considering that the annotators were not professionals.

## Methods and Evaluation Metrics

Supervised WSD was performed on the annotated corpora using Support Vector Machines (SVMs). SVMs were cho-

<sup>7</sup><http://www.elgalepin.com/>. The dictionary has been released online in 2007 and counts around 37.000 entries.

<sup>8</sup>We used the MultiWordNets on-line interface available at <http://multiwordnet.fbk.eu/english/home.php> (Artale 1997)

<sup>9</sup>The annotators were a 23 and a 54-year-old women, with no special linguistic training. Both are native speakers of Venetan.

Sense	%
<b>Omeni</b>	
1. Adult male person (opposed to woman)	14
2. A human being	59
3. Soldier	27
<b>Parte with coarse-grained labels</b>	
1. Something less than the whole	61
2. Role	8
3. Road or path (generic)	31
<b>Parte with fine-grained labels</b>	
1a. Something determined in relation to an entity	31
1b. Region or state	9
1c. One of the portions into which something is divided and which together constitute a whole	21
2. Role	8
3a. A line leading to a place or point	18
3b. Space for movement	13

Table 3: Sense labels and distribution

sen for a number of reasons: first of all, according to Navigli (2009), they achieve the best results in WSD compared to several other supervised methods. In particular, SVMs work efficiently in environments where there are a large number of features (Cabezas 2001) and are usually more resistant to overfitting (Lee 2004). Moreover, as stated in Yarowsky (2010), SVMs often perform well with few training examples per label. We used the implementation of LinearSVM in scikit-learn<sup>10</sup>. In the following subsections, the extracted features and the evaluation metrics are described.

### Feature Design

For designing features, we mainly followed Lee (2004) and Cabezas (2001). Overall, six features were implemented:

- Unordered bag-of-words (BoW) vector, considering all the lowercased words in the sentence.
- Unordered BoW vector with stopwords removed (as done by Lee 2004). We obtained a *stopword* list by selecting the most common tokens whose POS tag was in a restricted list (including articles, pronouns, clitics and conjunctions).
- Unordered BoW vector considering the wide-context (two sentences preceding and following the occurrence), which has been proved to improve noun disambiguation (Yarowsky 2010).
- Unordered BoW vector of the wide-context after filtering out *stopwords*.
- POS tags of the three words preceding and following the occurrence. The POS tag of the null token was denoted with a special symbol.
- Ordered sequence of tokens in the local, narrow context of the occurrence. Following Lee (2004), 11 features were developed, corresponding to different collocations.

<sup>10</sup><http://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html>

### Evaluation Metrics

For each feature combination, the following evaluation metrics were calculated using scikit-learn: Accuracy, Precision, Recall and balanced F1 Measure. These metrics were compared with upper and lower bounds. As baselines, two dummy classifiers were considered: the first one performs random classification, whereas the second one always chooses the most common class. Judges agreement rate was taken as the upper bound.

### Experiments and Results

We performed experiments using different feature combinations. As the amount of data available was scarce, a 10-fold cross validation strategy was used. In general, WSD tasks are difficult due to the very unbalanced sense distribution. Results with different feature combinations are reported in Table 4. Considering the word *omeni*, the most informative feature combination was BoW with the wide-context and the POS feature. As shown in Table 5, using this combination returns lower Recall for the first sense of *omeni*, but considerably improves all measures for the third sense, which was the lowest represented. In fact, occurrences of this sense usually occur in texts about American History, for which the wide-context feature can be useful in disambiguation. On the contrary, first sense occurrences often appear in quotations, which are unrelated to each other and for which the wide-context feature can be misleading.

Moving to the word *parte*, classification with coarse- or fine-grained labels follows the same pattern. In contrast to what happened with *omeni*, considering the wide context has a negative effect on the overall accuracy, whereas information about POS tags seems to be useful (see Table 4). However, best results are achieved using only local collocations together with the BoW feature. In fact, the corpus for the word *parte* is highly repetitive, so that considering ordered sequences of words near to the occurrence, like in the local collocations feature, can be very useful for disambiguation. As shown in Table 5, using this feature combination leads to acceptable Recall even for the second sense, which had very few occurrences.

	<i>omeni</i>	<i>parte I</i>	<i>parte II</i>
<i>Random classifier</i>	0.43	0.32	0.10
<i>Most common classifier</i>	0.60	0.64	0.32
BoW	0.77	0.81	0.58
Bow + wide_context	0.81	0.80	0.57
BoW+POS	0.73	0.85	0.62
BoW+wide_context+POS	<b>0.82</b>	0.85	0.62
BoW+collocations	0.75	<b>0.89</b>	<b>0.70</b>
BoW+collocations+POS	0.81	0.88	0.70
<i>Judges agreement</i>	0.89	0.93	0.79

Table 4: WSD Accuracy results using different feature combinations for *omeni*, *parte* with coarse labels (I) and *parte* with fine grained labels (II), compared with the baselines and upper bound.

	1st sense		2nd sense		3rd sense	
	P	R	P	R	P	R
<i>omeni</i>						
baseline	0.62	0.36	0.76	0.95	0.89	0.62
best comb	0.67	0.29	0.80	0.97	0.95	0.81
<i>parte I</i>						
baseline	0.86	0.92	0.00	0.00	0.76	0.82
best comb	0.89	0.97	0.71	0.56	0.94	0.82

Table 5: Comparison of Precision and Recall using only BoW feature and with the best combination. Results for *parte* with fine-grained labels follow the same pattern and are not reported due to lack of space.

## Conclusions and Future Work

In this paper, we reported on our annotation of a gold-standard and on the results of supervised WSD considering two Venetan words.

The annotation phase was laborious and time-consuming. As we were dealing with a LRL, we were not able to find expert annotators of Venetan, such as lexicographers (as in Wilks (1998)). Following Gale (1992), the difficulties of working with non-professional annotators were partially solved by redefining word-sense classification to a discrimination task. Concerning the results of WSD, we observed that different feature combinations performed better for a specific word. This is consistent with Resnik’s statement (Resnik 1997), according to which disambiguation, as a highly lexically sensitive task, in effect requires a specialized disambiguator for each considered word. In contrast with works on WSD for Old Languages and for other LRLs, (e.g., (Wunderlich 2015)), in our work we were able to access information from POS-tags. But, contrary to what we expected, considering POS-tags was not decisive for improving disambiguation. Furthermore, filtering out stop-words from the BoW features was not helpful for disambiguation. This experiment could be repeated using a different strategy to obtain a *stopwords* list. Additional knowledge sources could further improve accuracy if more tools for Venetan become available in the future. Future work could also investigate the adaptation of existing NLP tools for Italian to Venetan, as the two languages show a high degree of similarity (Tonelli 2010). It could be particularly interesting to obtain information about syntactic relations, which have been shown to be very discriminative in WSD (Lee 2004).

Overall, our proof-of-concept results are promising and demonstrate that, even with limited resources, the problem of WSD for an Italian dialect can be concretely approached, and we hope that our work will encourage further work on European dialects.

## Acknowledgments

The authors are thankful to Prof Rodolfo Delmonte of the Ca’ Foscari University of Venice for providing us with the STILVEN corpus.

## References

- Artale, A., Magnini, B., and Strapparava, C. 1997. Lexical discrimination with the Italian version of WordNet. In *ACL Workshop on Information Extraction*
- Berruto, G. 2005. Dialect/standard convergence, mixing, and models of language contact: the case of Italy. *Dialect change. Convergence and divergence in European languages*, 81-97.
- Bortolotti E., Rasom, S. 2005. The Ladin language between fragmentation and standardization: the contribution of Computational Linguistics, *Lesser Used Languages and Computer Linguistics, LULCL*
- Cabezas, C., Resnik, P., and Stevens, J. 2001. Supervised sense tagging using support vector machines. In *The Proceedings of the Second International Workshop on Evaluating Word Sense Disambiguation Systems: 59-62*. Association for Computational Linguistics
- Delmonte, R., Bristot, A., Tonelli, S., and Pianta, E. 2009. English/Veneto resource poor machine translation with STILVEN. In *Proceedings of the International Symposium on Data and Sense Mining, Machine Translation and Controlled Languages (ISMTCL)*
- Gale, W., Church, K., and Yarowsky, D. 1992. Estimating upper and lower bounds on the performance of word-sense disambiguation programs. In *Proceedings of the 30th annual meeting on Association for Computational Linguistics: 249-256*. Association for Computational Linguistics
- Jaber, S., Tonelli, S., and Delmonte, R. (2011). Venetan to English Machine Translation: Issues and Possible Solutions. In *8th International NLPCS Workshop: 69* ISO 690
- Lee, Y. K., Ng, H. T., and Chia, T. K. 2004. Supervised word sense disambiguation with support vector machines and multiple knowledge sources. In *Senseval-3: third international workshop on the evaluation of systems for the semantic analysis of text: 137-140*
- Navigli, R. 2009. Word sense disambiguation: A survey. *ACM Computing Surveys CSUR* 41(2), 10. ISO 690
- Resnik, P., and Yarowsky, D. 1997. A perspective on word sense disambiguation methods and their evaluation. In *Proceedings of the ACL SIGLEX workshop on tagging text with lexical semantics: Why, what, and how: 79-86*
- Singh, A. K. 2008. Natural Language Processing for Less Privileged Languages: Where do we come from? Where are we going? In *IJCNLP: 7-12*
- Tonelli, S., Pianta, E., Delmonte, R., and Brunelli, M. 2010. VenPro: A Morphological Analyzer for Venetan. In *LREC*
- Wilks, Y., and Stevenson, M. 1998. Word sense disambiguation using optimised combinations of knowledge sources. In *Proceedings of the 17th international conference on Computational linguistics (2):1398-1402*. Association for Computational Linguistics.
- Wunderlich, M., Fraser, A., and Langeslag, P. S. 2015. God Wat Thaet Ic Eom God - An Exploratory Investigation Into Word Sense Disambiguation in Old English. In *GSCL: 39-48*
- Yarowsky, D. 2010. Word sense disambiguation. In *Handbook of Natural Language Processing, Second Edition 315-338*. Chapman and Hall/CRC.