


What is natural language? Differences compared to artificial languages, and consequences for natural language processing

Diana Santos
Linguateca
www.linguateca.pt


 Information and Communication Technologies 1

What is natural language?

- Natural language is the oldest and most successful knowledge representation language
- Used for communication, negotiation, and reason (->logic)


What is natural language processing?

- Use computers to do things with natural language that should be useful for humans

 Information and Communication Technologies 2


How is natural language used?

- Most intelligent human tasks involve language
 - as center (communicating, teaching/learning, converting)
 - as periphery (mathematics papers, medical diagnosis, programming)
- Daily tasks
 - writing (and creating or conveying information or affection)
 - reading (and finding information)
 - translating (and mediating)

 Information and Communication Technologies 3

What are artificial languages?

- Languages consciously created by Man
- with the purpose of
 - reduce NL's expressivity
 - correct NL's problems
 - diminish NL's complexity
- by
 - abolishing or restricting ambiguity
 - giving a precise meaning to the vocabulary
 - a priori restricting the kind of communication
 - preventing change or evolution

 Information and Communication Technologies 4


What is a knowledge representation?

Five different roles:

1. a surrogate
2. a set of ontological commitments
3. a fragmentary theory of intelligent reasoning
4. a medium for efficient computation
5. a medium of human expression (!)


inevitable as long as we need to tell the machine about the world and as long as we do so by creating and communicating representations

Davis, Randall, Howard Shrobe & Peter Szolovits. "What is a knowledge representation?". *AI magazine* 14, 1, 1993, pp. 17-33.

 Information and Communication Technologies 5

Several natural languages

- There are different natural languages
 - reflecting a different world view
 - containing different "glue" (syntax, discourse)
 - taking into account different implicit information
- A theory of natural language evolution and how NL works ought to explain this undeniable fact
- Huge political impact of "taking care" of one's language

 Information and Communication Technologies 6

The rationale for *Linguateca*

- To do computational processing of Portuguese is not “just” do exactly the same applied to another language
- The lexicon is structured in a different way
- The grammar is different
- The culture(s) are different
- The technological shortcuts do not need to be the same
- One should apply reasoning and methodology appropriate to the object of study instead of uncritically assume that what has been used for English is appropriate for Portuguese

Santos, Diana. “Toward Language-specific Applications”, *Machine Translation* 14 (2), June 1999, pp. 83-112.

Linguateca, a project for Portuguese

- A distributed resource center for Portuguese language technology
- POSI project with FCCN as main contractor (2000-2006)
- First node at SINTEF ICT, Oslo, started in 2000 (work at SINTEF started 1998 as the *Computational Processing of Portuguese* project)

IRE model

- Information
- Resources
- Evaluation

www.linguateca.pt



Linguateca highlights, www.linguateca.pt

- > 1000 links More than 2,000,000 visits to the Web site
- [AC/DC](#), [CETEMPúblico](#), [COMPARA](#) ... Considerable resources for processing the Portuguese language
- *Morfolimpiadas* The first evaluation contest for Portuguese, followed by **CLEF** and **HAREM**

- Public resources
- Foster research and collaboration
- Formal measuring and comparison
- One language, many cultures
- Cooperation using the Internet
- Do not adapt applications from English

“General” view of natural language

- *Natural language contains too many ambiguities to be used alone for precise communication. Furthermore, natural language changes too rapidly for unqualified confidence in it as a preservation foundation.* (Gladney & Lorie, 2005:305).

Gladney, H. M. & R. A. Lorie. “Trustworthy 100-Year Digital Objects: Durable Encoding for When It’s Too Late to Ask”, *ACM Transactions on Information Systems* 23, 3, July 2005, pp. 299-324.

I will try to convince you of the opposite!

Presentation map

- Present some properties of natural languages
 - and their possible counterpart in other languages
- Touch upon the interrelationship of those properties
- Discuss the issue of many natural languages
- Tentatively conclude with some work that tries to deal with them

The list of properties

1. Metaphorical nature
2. Context dependency
3. Reference to implicit knowledge
4. Vagueness
5. Dynamic character (evolution and learnability)

1. Metaphor: constitutive property of NL

- Paradigm shift by Lakoff & Johnson (1980), *Metaphors We Live By*.
- *Most of our normal conceptual system is metaphorically structured, that is, most concepts are partially understood in terms of other concepts. (...) Understanding takes place in terms of entire domains and not in terms of isolated concepts.*
- Meaning is not a question of truth values (conditions for a sentence be true). It is necessary to understand first, in order to assign a truth value afterwards:
- *The true statements that we make are based on the way we categorize things and, therefore, on what is highlighted by the natural dimensions of the categories.*

Meaning and truth

- Truth depends on categorization in the following four ways:
 - a statement can only be true relative to some understanding of it.
 - understanding always involves **human categorization**
 - the truth of a statement is always **relative to the properties** that are highlighted by the categories used
 - categories are neither fixed nor uniform

Relativity's theory was right at hand for Heisenberg

Kinds of metaphors

- Conventional
 - structural
 - argumentation is war
 - orientational
 - more is up
 - ontological
 - love is a thing
- Imaginative or non-literal
 - new
 - extends old
 - uses unused part

Metaphor as culture

- Culture, our way of seeing the world, is embodied in metaphors
- Problems
 - as puzzles
 - as chemical precipitates
- Time
 - is money
 - is place

Artificial languages including metaphor?

- No, but understanding metaphor has already been useful for CS: Maglio, Paul P. & Teenie Matlock. "Metaphors We Surf the Web By", Paper presented at the *Workshop on Personalized and Social Navigation Information Space*, 1998, Stockholm, Sweden.
- spatializing UIs...
- *the terminology in [...] user interface design is acknowledging the central role of cognition in general and of metaphors in particular*
- Kuhn, W. "7+/-2 Questions and Answers about Metaphors for GIS User Interfaces", in Nyerges, T.L., et al. (eds.), *Cognitive Aspects of Human-Computer Interaction for Geographic Information Systems*, Kluwer Academic Publishers, 1995, pp. 113-122.

2. Context dependency

- The meaning of words, sentences and utterances **always** depends on context (which is sometimes the co-text)
- deictics, quasi-deictics, anaphors, etc.
- *aqui, ali, lá* (position), *eu, tu, nós* (speaker), *hoje, agora* (time)
- *depois, antes*
- *outro, este, primeiro*
- *presidente da República, presidente da mesa, o país, a assembleia,*
- *pai, meu*

o livro: article noun or pronoun verb?

- o livro da
 - o livro da cadeia
 - o livro da cadeia em
 - que o livro da cadeia em chamas
 - disse que o livro da cadeia em chamas
 - me disse que o livro da cadeia em chamas
 - me disse que o livro da cadeia em chamas só
 - me disse que o livro da cadeia em chamas só se
 - me disse que o livro da cadeia em chamas só se chamar
1. ...“Ivanhoe” é devido a... 2. ...a polícia!

Context dependency in programming languages

- Yes, also in artificial languages the “same” word means different things dependent on the context
- In $a=5$ and $b=a$, a in general means different things
- $a+3$ and $\text{substr}(a)$ and here also (coercion)
- $a=“aaa”$; $\text{print } a$ and $a=“bbb”$; $\text{print } a$, gives a different result
- just like print “print” or even $\text{print PRINT “print \$print”}$
- Dynamically typed languages and polymorphism (type-checking during evaluation), as well as ordinary scope rules.
- “Environments” in denotational semantics (symbol tables)

Context dependency 2

- Exactly the same sentence may mean the opposite
 - *Está calor aqui*
 - *Gosto imenso da comida portuguesa*
- Natural language presupposes an interlocutor
- There is a dialogue, a game of understanding something, explain better, change in the middle, leave things implicit
- What’s in a game: rules, cheating, negotiation, change, follow strategy, have a goal
- There’s external grounding...

Language is grounded in place and time

- *Every act of communication takes place in a material situation that plays an essential role in that communication.*
- *Much of what is now called context are really acts of communication*

Clark, Herbert H. “Pointing and placing”. In S. Kita (Ed.), *Pointing. Where language, culture, and cognition meet*. Hillsdale NJ: Erlbaum, 2003, pp. 243-268.

Context dependency 3

- Natural language is a language for **human** communication
 - *Although IT is a wonderful facilitator of data and information transmission and distribution, it can never substitute for the rich interactivity, communication and learning that is inherent in dialogue*
- Fahey & Prusak. “The eleven deadliest sins of knowledge management”, *California Management Review* 40 (3) 1998, 265-76.
- Conversations are the product of people engaging in joint activities
- Clark, Herbert. “Conversation, structure of”. In L. Nadel (Ed.) *Encyclopedia of Cognitive Science*. Basingstoke, England: Macmillan

Dialogue

- *Jet vet inte vad jag har sagt innan du har svarat och du vet inte vad du har sagt innan jeg har svarat. Du visar mig vad jag har sagt og jag visar dig vad du har sagt* (Molander 1996)
 - I don’t know what I have said until you answer and you don’t know what you have said until I answer. You show me what I said, and I show you what you said.
- Molander, Bengt. *Kunnskap i handling*. Göteborg: Daidalos, 1996.

Coercion: a way of modelling the role of context

- If you apply an operator, or something that expects different company, you may force a new behaviour: Treatments of tense and aspect use the concept of coercion to explain aspectual shifts
 - *functions which “coerce” their inputs to the appropriate type, by a loose analogy with type-coercion in programming languages (Ait-Kaci, 1984)*
- Moens, Marc & Mark Steedman. "Temporal Ontology and Temporal Reference", *Computational Linguistics*, Vol 14, Number 2, June 1988, pp. 15-28.

Context dependency implies non-compositionality

- Language as a set of building blocks each with a meaning and with a set of building rules which compute the meaning in terms of their parts (and the structure): **the compositional ideal**
- Language where all building blocks are interdependent and receive part of their meaning from the context, from the knowledge of the speaker and of the hearer etc.: **natural language**
- the ANY-thesis and the methodology of linguistics (Hintikka, 1980)
- *any* is grammatical when it has a meaning different from that of *some* in the context

An example of geographical IR

- Wishing to develop “non-geographic” topics restricted to a location in GeoCLEF, we found out that location was dependent on the topic
- “Europe” for UEFA and the international song contest
- former “Eastern bloc” may or may not include Germany
- and geography is dependent on time as well (Spain vs. Iberian Peninsula; India vs. Pakistan)
- while time is dependent on geography (Bronze Age, ...)

Everything is linked

- instead of ACTION in a PLACE in a TIME
- a more correct pseudo-formal way to put it would be **action(time,place) & place(time, action) & time(place, action)**
- And this is just to bring the traditional variables in, there are many other axes possible!
- Context is relevant to understanding of all parts

3. Reference to implicit knowledge

- You never start from nothing (except with babies)
 - An interaction is embedded or couched in a set of knowledge pieces
 - If one does not know anything (has no background) one cannot understand a text
 - Machine translation works thanks to implicit knowledge
- Schubert, K. "Implicitness as a guiding principle in machine translation", *Proceedings of COLING'88* (Budapest, 22-27 August 1988), pp.599-601.
- Translation is not possible without common knowledge
- Quine, Willard Van Orman. *Word and Object*, The MIT Press, 1960.

Extract explicit knowledge?

- *No matter how large our corpus, if it is domain specific, the major part of the domain ontology will not be specified because it is taken as given, or assumed to be part of the **background knowledge** the reader brings to the text.*
- Brewster, Ciravegna & Wilks. "Background and Foreground Knowledge in Dynamic Ontology Construction: Viewing Text as Knowledge Maintenance", in Ding et al. (eds.), *Semantic Web Workshop SIGIR 2003* (Toronto, July 28-August 1, 2003).
- *A text is an act of **knowledge maintenance**. (...) A primary purpose of a text at some level is to **change** the relationship between existing concepts, or **change** the instantiations of these concepts [...] or **adding** new concepts to the existing domain ontology*
- Brewster, Iria, Ciravegna & Wilks. "The Ontology: Chimaera or Pegasus", *Dagstuhl workshop on Learning for the Semantic Web*, 13-18 February 2005.

Separation from natural language

- How separate are actually artificial languages from NL? How much implicit information from actual English, Italian or French (Nirenburg & Wilks, Santos, etc.) is attached to:
- programming code
- mathematical formula
- KR languages
- musical notation, classical ballet

- specialized language texts: something in between?

4. Vagueness: the most important property

- The same unit means more than one related thing, at the same time.
 - Crucially different from ambiguity:
 - although both give more than one translation to one entity
 - the difference is in the **relationship** among the translations
 - vagueness is **systematic**, ambiguity is accidental
 - Vagueness has been the subject of much linguistic-philosophical research (Quine, Dahl, Lakoff, Kempson, Lyons, Keenan, etc. etc.) but it is somehow considered a nuisance for NLP
- Santos, Diana. "The relevance of vagueness for translation: Examples from English to Portuguese". *TradTerm* Vol. 5.1, 1998, pp. 41-98.

Examples of vagueness in Portuguese

- *João é amigo do Pedro* (Is *amigo* N or ADJ?)
- *Apaixonado, recusou o convite* (estar ou ser?)
- *Conhecer pessoas como ele é uma aventura* (conhecer bem, ou encontrar?)
- *Passando por casa dela, lembrei-me do irmão* (temporal ou causal?)
- *Encontraram-se na praia* (um ao outro, ou ambos lá?)
- *A porta abriu-se!* (sozinha, ou alguém a abriu?)
- *Ele quer casar com uma rapariga bonita* (que é, ou que seja?)
- *O homem que matou X é louco* (de re ou de facto?)
- *ou exclusivo ou inclusivo*

Vagueness and relationship with context

- *Ontem encontrei dois casais com o mesmo problema* (juntos ou separados?)
 - *Todos nesta sala sabem duas línguas* (as mesmas ou diferentes?)
- The context may select one interpretation, but may keep both
- *A construção levou muito tempo*
 - *A construção ficou muito bonita*
 - *A construção prejudicou-o imenso*

 - *Conheço o JJ há cinco anos*
 - *Conheci o JJ há cinco anos*
 - *Viver no Rossio há cinco anos deu-me o direito de participar.*

Kinds of ambiguity and vagueness

- **Syntactical** ambiguity or vagueness: more than one syntactical analysis or classification
 - **Semantic** ambiguity or vagueness: more than one sense/meaning
- Often both occur, but not necessarily:
- same syntactical analysis, semantically ambiguous: *sentei-me no banco.*
 - same semantic analysis, syntactically ambiguous: *sentei-me num banco na cozinha*
 - **Lexical** ambiguity or vagueness (a different thing: that is "located" in the lexicon)
 - **Contrastive** ambiguity or vagueness: more than one translation

Vagueness, polysemy and underspecification

- Vagueness is the general property: positively meaning related things
- Polysemy is vagueness restricted to the lexicon (related word senses)
- Underspecification is a more general name that includes vagueness: one might say that e.g. *table* is unspecified wrt weather, but not vague about the weather

- Vagueness is essential for communication, learning and evolution...

“Irreference”

The possibility of talking about non-existing things or intensions

- *O presidente dos Estados Unidos vai encontrar-se com o primeiro ministro do Canadá em Março de 2007*
- *Os leões são perigosos*
- *tenho comido muitos queques ultimamente*
- *eu como ao meio dia*

- *o fato de banho está algures no sótão*
- *quem quer que venha será recebido a tiro*
- *O João saiu antes de alguém chegar*

Vagueness in programming languages

- abstraction, encapsulation, lazy evaluation, loaded symbols
- polymorphism in Smalltalk -- you send the message *print*, and each receiver has its own version (implementation) of *print*: *A message specifies which operation is designed, but not how that operation should be carried out* (Goldberg & Robson, 1983)

Goldberg, Adele E. & David Robson. *Smalltalk-80: the language and its implementation*. Addison Wesley, 1983.

5. Language evolution

- origins (of underlying mechanisms)
- emergence (of specific features)
- evolution proper

Vogt, Paul, Bart de Boer & Tony Belpaeme. “Modelling language origins and evolution”, Tutorial, IJCAI 2005 (31 July 2005, Edinburgh, Scotland).

- evolution of communication
- grounding (relating language to the world)
- computational simulation, the emergence of compositionality

Examples of language evolution

- *viuvar* -> *enviuvar*
- *conheço* -> *reconheço*
- *costumar-se* -> *acostumar-se*

- *amar hei-de* -> *amarei*
- cases in Latin -> articles in Romance languages
- *amara* -> *tinha amado*

- *buscar* -> *ir buscar*
- *jogar a* -> *jogar à*
- *amar a* -> *amar*

Programming language evolution

- versioning
- deprecated syntax
- documentation
- compilers' warning messages
- upwards compatibility
- solving errors in previous versions

- different language paradigms (imperative, functional, ...)
- object orientation, extreme programming, ...

Difference between language designers and speakers

Acquisition and learnability of language

- In order to learn language, one has to use it: thus the primacy of dialogue and of context
- Every generation learns it anew...
- In order to learn, one has to actively extend it and see whether the extensions are sanctioned or not by future dialogues (Sampson's little Popperian in *Educating Eve*)
- Vagueness is a tool to help learning: not more than necessary needs to be decided/conveyed in an interaction, things can be refined, if needed, in dialogue

Representation and reasoning...

- ... are inextricably intertwined (Davis et al. 1993)
- Intelligent reasoning must be mirrored by natural language
- Appropriateness of inferences is more relevant than their validity
- *there is no reason to believe that systems for which notions like deductive closure are important have any demonstrable relationship to NLP, either as an empirical, engineering task or as a model of human processing*

Nirenburg, Sergei & Yorick Wilks. "What's in a symbol: Ontology, representation, and language". *Journal of Experimental and Theoretical Artificial Intelligence* 13 (1), pp. 9-23

The list of properties: again

1. Metaphorical nature
2. Context dependency
3. Reference to implicit knowledge
4. Vagueness
5. Dynamic character (evolution and learnability)

Connection among the several properties

- Vagueness and context dependence explain meaning shifts and therefore language evolution
- Dialogue explains how language can be learned, and how people can learn language by negotiating meaning
- Language builds upon implicit knowledge, and develops in the direction of the needs to communicate of a community; this explains divergence and convergence of languages
- Shared metaphors allow meaning shifts without trauma

Different natural languages: causes

- different contexts; different speakers
 - different implicit information; creative use every single minute
- Is it surprising that languages diverge?

■ *it would surely be surprising, and a very strong empirical claim, that different languages using different means to express 'meanings' always arrived at exactly the same end*

Keenan, Edward. "Some Logical Problems in Translation", in Guenther & Guenther-Reutter (eds.), *Meaning and Translation: Philosophical and Linguistic Approaches*, Duckworth, 1978, pp.157-89.

Different natural languages: results

- different grammar
- different lexical items
- different implicit/explicit conventions
- different rules for obligatoriness and optionality
- different discourse strategies
- different things we do with words
- different conventionalized metaphors
- different realities described/readily invoked
- different sayings/culture

Consequences to NLP

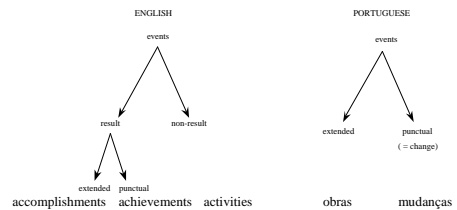
- Basically, it is essential to know how NL works in order to handle it
- Extreme simplification or denial of NL's properties won't do:
 - Deal with vagueness
 - Deal with metaphor
 - Deal with implicit knowledge
 - Deal with context
- To understand communication and dialogue is much more important than uncover the "right" syntax
- Many different factors work in tandem

Attempts to deal with vagueness

- In annotation, leave room for more than one category: HAREM and COMPARA
 - do not force a choice when it is not required
- Identify contrastively vague categories in tense and aspect
 - not only coercion
 - also aspectual classes or grammatical operators that can simultaneously mean more than one thing
- The translation network
 - linking two systems with different vague categories
 - explaining and formalizing concrete translation issues

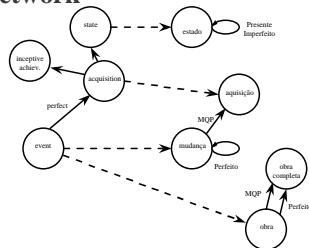
Categorization of events and states

- Different categories as far as events and states are concerned
- Different operators
- Attention to different details



The translation network

- Many ways to translate the pluperfect
- Transl. arcs are never meaning preserving



Santos, Diana. *Translation-based corpus studies: Contrasting English and Portuguese tense and aspect systems*. Amsterdam/New York, NY: Rodopi, 2004.

Let us look at the translation *sit* – *sentar*

45 translationally relevant different senses of *sit* in COMPARA

- sit up, sit through, sit with, sat on, sat forward, sitting up, sit out
- sit in judgement, sit on the bench
- sit to
- there sit, still sit
- sit things out, who should sit where, sits on the skyline, sat deep in discussion, sat down to a meal, sitting round, sitting about all over...
- sat Xing, sat opposite, behind, next, alone, attentively, bent...

70 different translations in COMPARA

What gets translated by *sit*?

123 different expressions/contexts get translated by *sit*

- ficar *ando
- mesa
- ficar
- ter à frente, ter diante
- de onde estava

many of them with absolutely no evidence for a sitting position

- fiquei só de cueca, gostava de olhar, fiquei com eles à mostra, fodendo num motel, velava o corpo, ali o deixaram ficar, e é a calma personificada, as pessoas não trazem muita pressa, benigna e sossegadamente determinemos, deixa-te estar af, folheávamos, a vossos pés...

What can be concluded?

- The more language one sees, the larger number of “senses” and or contexts relevant for understanding/translating. In a Zipfian way? Kilgarriff, Adam. “How dominant is the commonest sense of a word?” In Sojka, Kopecek & Pala (eds.), *Text, Speech, Dialogue*. Springer, 2004, pp.103-12.
- Fixed inventory of senses (or translations) is never the whole story!
- Huge amount of context and implicit information that also plays a role, together with language habits and culture (metaphor, idioms, ...)

Conclusion

- Natural language is fundamentally different from artificial languages, so far
- It has properties that developed and were refined through the whole history of mankind
- If we want computers to be efficient and helpful mediators for the human race, natural language processing has to deal with language and understand it
- It is possible that by studying NL we are also able to devise better artificial languages or improve reasoning capabilities of our current systems

Three kinds of vagueness

1. monodimensional, fuzzy borders
 - how many hairs can you have and still be called bold?
2. set of properties in different axes
 - Swedish citizen: Swedish parents and born in Sweden
 - noun: inherent gender; no degree; countable, mass or abstract
3. set of related properties
 - acquisition class (to mean the state and its inception)
 - privative opposition
 - gosto: bom gosto, mau gosto
 - sorte: boa sorte, má sorte