# SEMANTIC CONSTRUCTION OF UNIVOCAL LANGUAGE

## Alejandro De Santos, Pedro G. Guillén, Eduardo Villa, Francisco Serradilla.

**Abstract**: In this paper a solution is propose to organize the space of words that exist in a specific language in their different semantic categories. By taking a natural language, we're going to define a unique meaning for each word, as a construction made (d,C) of pairs of words and contexts. On the other hand, let us consider the space of meanings. All the words that share meaning (synonymous words) can be associated with one meaning. This permitus to make a partition of the space of words in groups of synonyms. Finally, a classification of the space of words will be obtained in the different groups of words that share meaning. This allow to choose the useful word that represent a meaning, and reduce the number of words selecting one representative from each group of synonyms. It will be very useful for calculating distances between words.

**Keywords**: Natural Language, Semantic, Context.

**ACM Classification Keywords**: C.2.1 Network Architecture and Design; D.2.1 Requirements Specifications (D.3.1) Languages; D.3.2 Language Classifications Applicative (functional) languages; E.1 [Data Structures] Graphs and networks; H.5.2 [User Interfaces] Natural language.

## Introduction

Let $A'$ be a natural language. We built $D = (A', *)$, with a free semigroup structure, treated as a set, regardless of its algebraic properties [Dieter, 2004]. Can be built the pair $(d, C)$ made of one word and one context that select one of the possible meanings of the word. So, a solution to the problem of poly semy is proposed. Now, is necessary to consider the space $\Delta$ of all the possible meanings that built a language. Over the pair $(d, C)$ the meaning mapping $\psi: (d, C) \longrightarrow \Delta$ is defined (we're going to include programming abstractions such as procedures, functions [Miguel, 2011]) $\psi$ assigns to each pair $(d, C)$ only one meaning. Finally, an equivalency relation over $(d, C)$ is defined in order to build a quotient space inducing a classification in the different semantic classes [Ito, 1977], [Angelova, 1988], wherein each equivalence class is formed by all the words associated with a determinate meaning (synonymous words) [Ito, 1981].

## Space of Words

***Definition* 1**. Let $\mathbf{A} = \{a, b, c \dots\}$ be an alphabet, such that $A \neq \emptyset$.

We can define a word as a finite succession of elements of $\boldsymbol{A}$, where repetitions are allowed.

For example **aaa**, **aabcc**, **b** are different words.

Let $A'$ be the set composed by words.

***Definition* 2**. Let $\boldsymbol{D} = (\boldsymbol{A'}, *)$ be a free semigroup with a associative law (*) juxtaposition, and an identity element $1 \in G$.

We can multiply two words **aaabbb*abccc=aaabbbabccc**, where the identity element '1' is defined as the empty succession of words.

The grammatical rules are defined as restrictions on this free semigroup.

**_Definition_ 3**. Let the inclusion mapping be $i_{|\Theta}: D \longrightarrow \Theta$ , $\Theta$ the space of concepts which belongs to a given lexicon $\mathcal{L}$. The inclusion mapping $i_{|\Theta}$ saves the words that belong to $\mathcal{L}$.

Let $d \in D$ then $i_{|\Theta}(d) = \{d/d \in \Theta\}$ and $i_{|\Theta}(d) = \{\emptyset/d \in D \setminus \Theta\}$.

For example $i_{|\Theta}(\boldsymbol{aabbacc}) = \emptyset$ such that $\boldsymbol{aabbacc} \notin \mathcal{L}$ and $i_{|\Theta}(\boldsymbol{dog}) = \boldsymbol{dog}$ such that $\boldsymbol{dog} \in \mathcal{L}$.

In the next point the space of context is building as a simple generation rule so as we will see how to calculate the shortest path between two contexts.

## Space of Contexts

Let $C$ be a context, a list of words that define an ambient where to locate a word belonging to $\mathcal{L}$.

**_Definition_ 4**. Let $C$ a graph made of $\{C_1, C_2, \dots C_n\}$ a countable set of contexts, such that $C_i \in \Omega$ and $\Omega$ is the space of contexts inspired by [Dieter, 2004].

We're going to build $C$ with a simple **_formation rule_**. Let $C_i \in \Omega$ be a determined context composed by more specific disjoint contexts $C_{i+1} \in \Omega$.

$$C_i = \bigcup_{i=1}^{n} C_{i+1}^j \ \ such \ as \cap_{j=1}^{n} C_{i+1}^j = \emptyset \equiv [4.1]$$

**_Definition_ 5**.Let the **_generation rule_** be the process of division of a context $C_i \in \Omega$ in their constituents subcontexts $C_{i+1} \in \Omega$ in [4.1]. The number of generation rules must be countable to avoid problems of computability.

We consider $C_i$ as a vertex of the graph and all the contexts in which $C_i$ is subdivided as a set $\{C_{i+1}^1, C_{i+1}^2, \dots C_{i+1}^n\}$ of new vertex's connected with $C_i$ by edges.

If we are applying the rule generation to $C_i$ and one of the subcontext generated $C_{i+1}^3$ is the same as another $C_{j+2}^5$ belonging to another division $C_{j+2}^5 = C_{i+1}^3$ , then this two subcontexts are considered the same and we can see it in the graph using the edges from the previous settings $C_i, C_{j+1}$ to it.

This breaks the tree structure of the graph, creating closed paths and cycles.

Let $\|\alpha(C_i, C_j)\|$ be the **_lenght of a path_** between two contexts [Dieter, 2004].

**_Definition_ 6**. The **_shortest lenght of a path between two contexts_** is:

$$\|\alpha_{min}(C_i, C_j)\| = \{\|\alpha(C_i, C_j)\| \ such \ that \forall \ \|\alpha'(C_i, C_j)\|, \|\alpha'(C_i, C_j)\| \geq \|\alpha(C_i, C_j)\|\}$$

It will be useful to calculate in the future, the minimum distance between two contexts.

In the next point we're going to build the relationship between the space of words $\Theta$ and space of contexts $\Omega$ in the pair $(d, C) \in (\Theta, \Omega)$.

Once this is achieved, we will assign one meaning to each pair through a well defined mapping.

## The Meaning Mapping

**_Lemma_ 1**. Let $(d, C) \in (\Theta, \Omega)$ be the pair made of one word and one context. $\Theta$ is the space of concepts which belongs to a given lexicon, and $\Omega$ is the space of contexts.

It is important to realize that most of the words of a given lexicon $\mathcal{L}$ are polysemous.

Is selected one of the possible meanings of the word through the context where to locate the word.

Let $\Delta$ be meanings space of a given lexicon $\mathcal{L}$.

***Lemma* 2.** Let $\psi: (\Theta, \Omega) \longrightarrow \Delta$ a mapping. $\Delta$ is the meanings space, $\Theta$ is the space of words that belong to $\mathcal{L}$, and $\Omega$ the space of contexts.

Once we have selected a word and a context by the pair, we assign one meaning to the pair through the mapping $\psi$.

For example, the word "paint", choose different meanings for different contexts through the mapping $\psi$.

$\psi(paint, decoration)$ =a coloured substance which is spread over a surface and dries to leave a thin decorative or protective coating.

$\psi(paint, basketball)$ =a rectangular area marked near the basket at each end of a court.

***Proposition* 1.** $\psi$ is a well-defined mapping.

***Proof*.** By making $\psi$ a mapping, we have found a way to avoid the polysemy problem and assign one only meaning to each pair (**word, context**) [Kazimierz, 2010].

Let $d \in \Theta$ be a word.

Let $C_1, C_2, \dots C_n$ be a list of contexts associated to the word $d$.

$$\forall (d, C_i) \, \exists! \, h \in \Delta / \psi(d, C_i) = h \text{ with } i \in \{1, 2, \dots n\}.$$

We can consider a pair $(d, C) \in (\Theta, \Omega)$, the support of $\psi$, that allow to define $\psi$ as a mapping.

$\psi: (\Theta, \Omega) \longrightarrow \Delta$ defines only one meaning for each pair.

In the last point we will organize the pairs $(d, C) \in (\Theta, \Omega)$ by groups that share meaning (synonymous words) through an equivalence relation $\sim$. This allow to choose the word that represents a meaning that suits us, and reduce the number of words by choosing one representative from each group of synonyms. It will also be very useful for calculating distances between words.

## The space of semantic meanings

***Lemma* 3.** Let $\sim$ be an equivalence relation. Two words in one of their concrete contexts are related: $(d, C) \sim (d', C')$ if and only if $\psi(d, C) = \psi(d', C')$.

For example:

$\psi(late, time) =$ after the expected or usual time.

$\psi(delayed, person) =$ retarded, incapacitated.

$\psi(delayed, traffic) =$ after the expected or usual time.

Where $\psi(late, time) = \psi(delayed, traffic)$.

***Teorem* 1.** $\sim$ is a **equivalence relation**.

***Proof*.**

$\sim$ is a **reflexive relation**: One word in one of its meanings is related with himself

$(d, C) \sim (d', C') \leftrightarrow \psi(d, C) = \psi(d', C')$

Obviously: $(late, time) \sim (late, time) \leftrightarrow \psi(late, time) = \psi(late, time)$ =after the expected or usual time.

$\sim$ is a **symmetric relation**: If one word, in one of their meanings, is related with another one (in one of its meanings). This involves that the another word is related with the first word, each one, in one of its meanings.

$(d, C) \sim (d', C') \leftrightarrow \psi(d, C) = \psi(d', C') \equiv \psi(d', C') = \psi(d, C) \leftrightarrow (d', C') \sim (d, C)$

Is not difficult to see:

$(late, time) \sim (delayed, traffic) \leftrightarrow \psi(late, time) \sim \psi(delayed, traffic) =$ after the expected or usual time$= \psi(delayed, traffic) = \psi(late, time) \leftrightarrow (delayed, traffic) \sim (late, time)$.

$\sim$ is a **transitive relation**: When one word(always in one of their meanings) is related with another and the third word is related with another one, this implies that the firs and the last word are related.

$\qquad (d, C) \sim (d', C')$ and $(d', C') \sim (d'', C'') \leftrightarrow \psi(d, C) = \psi(d', C') = \psi(d'', C'')$

In this case:

$(late, time) \sim (delayed, traffic)$

and $\qquad (delayed, traffic) \sim (overdue, born) \leftrightarrow \psi(late, time) = \psi(delayed, traffic) =$ $\psi(overdue, born)$ that means as we know: after the expected or usual time.

This **equivalence relation** $\sim$ allow to organize the different pairs in their classes:

$[(d, C)] = \{(d', C') \in (\Theta, \Omega)/(d', C') \sim (d, C)\} \equiv \{(d', C') \in (\Theta \times \Omega)/\psi(d', C') = \psi(d, C)\}$

consisting of all possible pairs, formed by words and their contexts that share the same meaning.

This **equivalence relation** $\sim$ induces in $(\Theta \times \Omega)$ the classification of their elements in the different semantic classes, building the quotient space $(\Theta, \Omega)/\sim = \{[(d, C)]/(d, C) \in (\Theta \times \Omega)\}$ wherein each element is one of the classes that share meaning.

For    example:

$[(late, time)] = \{(delayed, traffic) \in (\Theta, \Omega)/(late, time) \sim (delayed, traffic)\}$.

We can see that:

$$\cup_{\forall\, i,j} [(d_i, C_j)] = (\Theta \times \Omega)/\sim \text{ such that } j \in \{1,2 \dots n\}, i \in \{1,2 \dots m\}$$

$$\cap_{\forall\, i,j} [(d_i, C_j)] = \emptyset \text{ such that } j \in \{1,2 \dots n\}, i \in \{1,2 \dots m\}$$

**Lemma4.** Let $\psi *$ be the meaning mapping defined on the quotient space $\psi *: (\Theta, \Omega)/\sim \longrightarrow \Delta$

$\psi *$ is a mapping from the classes to the meanings

$\qquad \psi * [(late, time)] = \psi(late, time) = \psi(delayed, traffic) = \psi(overdue, born)$

that means as we know: after the expected or usual time.

**Theorem2.** The equivalent relation $\sim$ makes be $\psi *$ a injective mapping.

**Proof.** $\psi * [(d, C)] = \psi * [(d', C')] \leftrightarrow [(d, C)] = [(d', C')] \equiv (d, C) \sim (d', C')$

Two classes have the same meaning if they are the same class.

Finally the injective mapping $\psi$ allows us to choose only one word that represents a meaning that suits us from each group of synonyms, simplifying the process.

We can see that:

$\psi * [(late, time)] = \psi * (delayed, traffic) =$ after the expected or usual time.

## Conclusion

We approach a solution to the problem of polysemy, building the different pairs $(d, C)$, choosing an unique meaning for each word. Acting on this pairs $\psi(d, C)$, the meaning mapping $\psi$, assigns one meaning to each

pair. Finally, we have solved the synonymy problem, organizing all the words in their different semantic classes, through the quotient $(\Theta, \Omega)/\sim$, where two pairs $(d, C) \sim (d', C')$ are related if and only if they share the meaning $\psi(d, C) = \psi(d', C')$. With this action, we reduce the number of words by choosing one representative from each group of synonyms, allowing the selection of the word that represents the useful meaning.
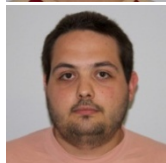
## Acknowledgements

## Bibliography

[Angelova, 1988] Angelova, GalyaMladenova. Syntactic and semantic structures for relational query languages . (Russian) Tanulmányok-MTA Számitástech. Automat.Kutató Int. Budapest No. 209 (1988), 187 pp. 68P15 (68Q55)

[Dieter,2004 ] Dieter Jungnickel. Graphs, Networks and Algorithms. Volume5. Ed. Springer, 2004

[Ito, 1977] Ito, Tetsuro; Kizawa, Makoto. Semantic structure of naturallanguage. Systems-Computers-Controls 7 (1976), no. 2, 110 (1977).

[Ito, 1981] Ito, Tetsuro; Toyoda, Junichi; Kizawa, Makoto, Hierarchical data base organization for document information retrieval.Systems-Comput.-Controls 10 (1979), no. 2, 39–47 (1981).

[Kazimierz,2010] Kazimierz Subieta. Syntax and Semantics of the Stack Based Query Language (SBQL)1, Polish-Japanese.Institute of Computer Science Polish Academy of Sciences Version 2, 18 June 2010. Page 41. Store mode Classes and inheritance.

[Miguel, 2011] Miguel Gonzalez. ESTRUCTURAS ALGEBRAICAS. 4de febrero de 2011 .Grupos libres. Generadores y relaciones. Página 93
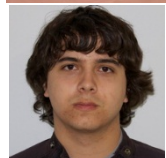
## Authors' Information

***Alejandro De Santos*** *– Natural Computing Group, Facultad de Informática, Universidad Politécnica de Madrid, Campus de Montegancedo s.n., 28660 Boadilla del Monte, Madrid, Spain; e-mail: matematicofeliz@gmail.com*



***Pedro G. Guillén*** *– Natural Computing Group, Facultad de Informática, Universidad Politécnica de Madrid, Campus de Montegancedo s.n., 28660 Boadilla del Monte, Madrid, Spain; e-mail: pedrogguillen@gmail.com*



***Eduardo Villa*** *– Natural Computing Group, Facultad de Informática, Universidad Politécnica de Madrid, Campus de Montegancedo s.n., 28660 Boadilla del Monte, Madrid, Spain; e-mail: e.villa.valdes@gmail.com*



***Francisco Serradilla*** *–MercatorGroup, Universidad Politécnica de Madrid Carretera de Valencia Km. 7, 28031 Madrid, Spain; e-mail: fserra@eui.es*