

# Similitud Semántica: Comparación y Crítica a los Modelos Actuales<sup>1</sup>

**Enrique P. Latorres**

Universidad ORT del Uruguay  
Montevideo, URUGUAY, 11000  
latorres@adinet.com.uy

## **Abstract.**

There are several reasons for considering semantic similarity as one of the most important problems of Information Technology today. Most tasks of knowledge evaluation and scientific research, or even human common sense knowledge are accomplished by some kind of semantic similarity matching. When talking about complex problem solving, most of its complexity is that of identifying the problem itself. This means that a conceptual understanding of the problem is necessary to match its specification to the expected “input” and conditions of the solving procedure. Future systems based on knowledge must be able to reuse knowledge from different sources and should handle aspects not considered in current models. In this document many models are analyzed together with critics to their implementations or theoretic justifications, and suggests attributes that should be included in such a new paradigm.

**Keywords:** Semantic Similarity, Agent Conceptualization, Artificial Reasoning, Restrictions, Semantic Reuse, Knowledge Integration.

## **Resumen.**

Hay varias razones para considerar el problema de Similitud Semántica como uno de los más importantes para la Tecnología de la Información. La mayor parte de las tareas de evaluación de conocimiento e investigación científica, o aún la aplicación de conocimiento de sentido común humano, son desarrolladas mediante algún tipo de mapeo de similitud semántica. Cuando hablamos de resolución de problemas complejos la mayor parte de la complejidad es la de identificar el problema. Esto significa que es necesaria una comprensión conceptual del problema para vincular su especificación con los parámetros y condiciones esperados para un mecanismo de resolución de problemas. En este documento se analizan muchos de los modelos utilizados actualmente junto con las críticas a sus implementaciones y justificación teórica, y sugieren atributos que este nuevo paradigma debería incluir.

**Palabras claves:** Similitud Semántica, Conceptualización por agentes, Razonamiento Artificial, Reutilización Semántica, Integración de Conocimiento.

---

<sup>1</sup> Este trabajo está parcialmente financiado por el Programa de Desarrollo Tecnológico (PDT) contrato S/C/BE/06/04. <http://www.pdt.gub.uy>.

## 1 Introducción

Hay varias razones para considerar el problema de Similitud Semántica como uno de los más importantes para la Tecnología de la Información. La mayor parte de las tareas de evaluación de conocimiento e investigación científica, o aún la aplicación de conocimiento de sentido común humano son desarrolladas mediante algún tipo de mapeo de similitud semántica [1][2], también conocido como “conceptual matching”. Cuando hablamos de resolución de problemas complejos, la mayor parte de la complejidad es la de identificar el problema. Esto significa que se necesita una comprensión conceptual del problema para vincularlo a la mejor forma de solucionarlo. Más directamente, para resolver un problema debemos obtener la mejor correlación entre la especificación del problema y los parámetros y condiciones esperadas para un mecanismo de resolución de problemas.

Uno de los problemas más complejos es el de reutilización de conocimiento o software [3], sobre lo que se ha trabajado mucho en variadas áreas. La idea de reutilizar el esfuerzo realizado en el pasado para disminuir los costos de los desarrollos futuros, ha estado en el tapete desde siempre. La puesta en práctica de esta idea es tan variada como actividades y paradigmas hay en el desarrollo de software. Se ha intentado reutilizar, con variada fortuna, subrutinas -en bibliotecas-, objetos, componentes, programas. Pero no solo el código se reutiliza, también el diseño (e.g. Los patrones o *patterns*) y también los requerimientos (e.g. La especificación de seguridad informática en el *Common Criteria*).

Las ideas de “qué es lo que debe reutilizarse” entre distintos proyectos de software es aún más un arte que una ciencia. Sin embargo algunos patrones parecen mantenerse:

- a) Si la reutilización es factible, la ganancia de rentabilidad crece con el nivel de abstracción del objeto usado.
- b) La factibilidad de la reutilización depende del “lenguaje” en el que se expresa lo que se espera reusar.
- c) La reutilización aparece tanto en cosas de audiencia amplia (e.g. en cosas comunes a muchos sistemas) como en cosas de dominio restringido (e.g. los programas de un cierto dominio de negocios).
- d) La reutilización depende de otras cosas que aún no sabemos prever.

En general cualquier mecanismo de reutilización implica un cierto nivel de clasificación y de similitud semántica a efectos de agrupar conceptos y asociar problemas a mecanismos de solución [4]. Esto nos lleva a buscar soluciones al problema de similitud semántica. Similitud Semántica es encontrar qué tan similar es un concepto, o una pieza de software, o de información, o servicio, o proceso a alguna otra, a efectos de poder integrar, interoperar y/o reutilizar, economizando en esfuerzo y recursos.

Estos procedimientos de reutilización y correlación problema-solución están basados principalmente en actividades y decisiones humanas. Por lo tanto son caras, lentas y con propensión a errores, y hace que el esfuerzo necesario para reutilizar no sea tan atractivo como debiera, debido al costo de implementar los procedimientos que garanticen la calidad de las unidades reutilizadas.

El problema de similitud semántica debiera ser uno de los principales para la ingeniería del conocimiento. Una solución adecuada para este problema podría permitir respuestas a muchos otros. Esto tiene que ver con resolución inteligente y automática de problemas, como ser integración, interoperación y reutilización automática de conocimiento y bases de datos, traducción automática de especificaciones formales en código, ingeniería reversa automática, reconocimiento de patrones, identificación y recuperación de información automática, razonamiento basado en casos, aplicabilidad de problema/solución, y muchos otros que por ahora necesitan supervisión humana. Es seguro que soluciones a este problema podrían cambiar la efectividad y la aplicabilidad de los sistemas de inteligencia artificial en el futuro[5].

Este trabajo es parte de una investigación en la búsqueda de modelos para resolver la similitud semántica en forma computacional y automática. Muchos de los conceptos aquí vertidos son un resumen de los resultados de esta investigación de tesis doctoral.

## 2 Similitud semántica

La investigación en similitud semántica tiene una larga historia pero solo recientemente fue estudiada en forma más o menos amplia. Los investigadores que trabajan en similitud cubren varias disciplinas, desde la psicología cognitiva, pasando por la neurofisiología, la filosofía y la lingüística, hasta la inteligencia artificial por nombrar algunas. Asimismo, las aproximaciones de las diferentes disciplinas tienden a influenciarse unas a otras por medio de referencias cruzadas en publicaciones e investigación. Por lo que diferentes ideas tienden a ser re-expresadas y aparecen en diferentes formas y relevancia.

Un iniciador de los modelos de similitud semántica fue Russell [6], con su teoría de clases y similitud. Desde el comienzo su intención de hacer un modelo que pudiera considerar todos los aspectos del mundo real fue en vano, al

ser sensible a varias paradojas. Una de estas paradojas, el problema de “la clase de clases”, lo compartió con Frege para encontrar una solución “Ninguna clase puede ser una extensión de una super-clase que cae debajo de la primera”, resolviendo el problema, pero haciendo que la solución sea mucho más restrictiva, y al final menos capaz de representar la realidad. Russell, trabajando por su cuenta, desarrolló la teoría de proposiciones, pero de todas maneras volvía a caer nuevamente en ciertas situaciones paradójicas. De todas formas, estos modelos afectaron la forma de percibir los problemas de la realidad entre los futuros filósofos, los lógicos, y desde ellos hasta los teóricos de la computación de hoy.

Price [7] agrega a este modelo el concepto de atributos que pueden compartir similitudes, y separó la visión ideal de especificaciones de la realidad con la visión parcial percibida por los humanos. Este autor concluye que la determinación de pertenencia a una familia por la experiencia, muchas veces es imposible, pero es alcanzada a través de la similitud de los atributos percibidos. Más aún, similitud puede no estar relacionada a una jerarquía de clases como ser el caso de una estatua y la persona que esta representa.

Osgood [8] analizó el significado de las palabras experimentando con la interpretación humana desde un punto de vista psicolinguista. Este tipo de evaluación, basado en el análisis de humanos realizando evaluación de sinónimos, también es realizado por Rubinstein y Goodenough [9]. Este trabajo fue validado luego por Miller y Charles [10]. Los resultados de estos experimentos son utilizados en varios estudios recientes para probar y validar estadísticamente varios algoritmos basados en computadora a efectos de determinar la similitud semántica entre conceptos expresados por palabras. Se fundamentan en la idea de que el juicio humano se considera correcto por definición, tal como es afirmado por varios autores [24]. Estos aseveran que la similitud es lo que hacen los humanos, por lo que comparan el desempeño de los algoritmos de computadora con el desempeño medido en una muestra de humanos realizando la tarea.

Si bien varios de estos autores expresan que hay algún tipo de relación inversa entre la similitud semántica y el concepto de distancia semántica, la fórmula exacta es muy discutida, y varía ampliamente de un autor a otro. Algunos de los modelos implementados pueden ser clasificados en pocas categorías: Taxonómicas, por Atributos, Mixtas, de Grafos de Dependencias, etc. Varias particularidades están presentes en varios modelos, por lo que es difícil hacer una clara separación entre ellos. La clasificación fue realizada principalmente basada en el estereotipo y sus atributos más relevantes.

Nótese que estos modelos de similitud devuelven valores que son diferentes del caso bivaluado:  $similar = \{verdadero, falso\}$ . En muchos casos se puede asociar un concepto *fuzzy* de similitud.

## 2.1 Modelos Taxonómicos

Muchos investigadores evalúan la similitud semántica de acuerdo a una jerarquía de clases preestablecida (taxonomía), algún tipo de ontología que especifica las relaciones entre conceptos. Las implementaciones de similitud semántica en el enfoque taxonómico involucran en general la cuenta de bordes o vértices (relaciones) entre dos conceptos con o sin algún tipo de ponderación en los enlaces relevantes y en otros atributos.

### 2.1.1 Cuenta de vértices (*Edge count*), Densidad Normalizada (*Normalized Density*), Densidad de Distancia (*Distance Density*)

Varios investigadores han basado sus modelos en la distancia entre dos conceptos mediante la cuenta de vértices topológicos que son recorridos en una taxonomía o red semántica, desde un concepto hasta otro, sobre la base de la sugerencia de Rada [11], generalmente con alguna ponderación afectando a cada vértice. Pequeñas modificaciones a este modelo tienen variaciones significativas en efectividad. También, la taxonomía y el tipo de relaciones consideradas afecta en forma importante la efectividad del modelo, como se puede ver en el trabajo de McHale [12].

Ramon y Van Laer [13] definieron una medida de distancia entre átomos de primer orden que tiene todas las propiedades de una métrica. En forma similar a otros modelos basados en taxonomías, la distancia es calculada al comparar la distancia de vértices a una superclase común con ponderaciones. Otro modelo es el de Distance-Density de Krumhansl [14]. En este se asume que en regiones densas de estímulo, o sea clases que tienen muchas subclases, se deben hacer discriminaciones más detalladas que en zonas menos densas.

### 2.1.2 Contenido Probabilístico o de Información

Resnik [15] sugirió que cada cuenta de vértices sea corregida por una ponderación probabilística estimada empíricamente. Este autor consideró que la información contenida en un mensaje se puede medir como el negativo del logaritmo de la probabilidad de que ocurra el mensaje, siguiendo la propuesta de Shannon [16], aunque no lo citó explícitamente. El contenido de información es calculado a partir de un objeto al analizar la probabilidad de pertenecer a una cierta super-clase. Esto puede ser visto también como una adaptación de conocimiento estadístico a las dependencias de las relaciones. Asimismo se puede considerar que hay un cierto nivel de razonamiento

(estadístico en este caso) en la evaluación de similitud. El modelo fue comparado con los resultados de Miller y Charles, los cuales tienen una correlación del 79%.

Resnik mostró algunos aspectos importantes sobre distancia semántica. Similitud Semántica es un caso especial de relación semántica, como en el caso de *car-gasoline* (automóvil-gasolina) con relación a *car-bicycle* (automóvil-bicicleta) sobre la base de Wornet[17]. Intuitivamente, los primeros dos conceptos están más relacionados pero el segundo par tiene más similitud. Ante esto plantea que no alcanza con considerar solo las relaciones de hiponimia e hipernimia (subclase y superclase) sino también la intención y uso del concepto, ya que su modelo se basa en la sugerencia de Rada [11]. El trabajo de Resnik es similar al de Leacock y Chodorow [18]. Estos últimos consideran el largo del camino normalizado al dividir por la altura máxima de la taxonomía. Resnik mostró como las conexiones entre conceptos tienen conceptualmente e intuitivamente una distancia diferente. También Leacock y Chodorow sugieren que [18], dependiendo del caso, cualquiera de ambos métodos puede tener mejor resultado, por lo que la conveniencia de usar el método de contenido de información depende de la taxonomía y la cantidad de datos para superar a la cuenta de vértices normalizada.

Dekang Lin [19] definió otro modelo de similitud normalizado basado en teoría de la información. Lin asume que el método debe derivarse de ciertas suposiciones sobre propiedades de la medida de similitud. Este modelo toma de Tversky [29] que en la medida de similitud debe considerarse tanto los atributos comunes como los que los diferencian. Lin mostró como su modelo es un *framework* para desarrollar otras medidas de similitud para diferentes tipos de elementos, incluyendo valores a los que se pueden asociar aseveraciones de cualidad difusa. El experimento en similitud de palabras brinda resultados impresionantes luego de analizar un corpus extenso basado en tripletes de dependencias. Esto mostró una gran potencia para la recuperación de corpus de grandes textos. Lin comparó sus resultados con los de Miller y Charles [20], y obtuvo una mejor correlación que otros modelos, tales como los de Resnik [15], Wu y Palmer [21], además de otras variantes de similitud como ser la distancia de cadenas de caracteres de Levenshtein [22].

Jiang y Conrath [23] utilizaron un modelo similar al de Lin, pero su pequeña modificación tiene resultados importantes de mejora. Curiosamente su modelo también recuerda a la fórmula de Shanon de transmisión de mensajes en un medio con ruido.

Budanitsky y Hirst [24] mostraron como esta simple modificación incrementa notablemente la efectividad de la medida de similitud. El modelo de Jiang y Conrath [23] presenta una correlación del 85% con el experimento de Miller y Charles. Esto se aproxima bastante bien al máximo estimado de 88% en la duplicación del experimento de Miller y Charles[10] realizado por Resnik [15] y al máximo estimado de 90%. Pero todos estos modelos son comparados contra experimentos basados en un muy pequeño corpus fijo.

### 2.1.3 Intencional

El modelo de Hakimpour y Geppert [25] está basado en la detección de relaciones de similitud o diferencia. Estas son descritas por relaciones de semántica intencional, proposiciones basadas en el modelo de Goh et al. [26]. Una sugerencia interesante en este trabajo es que se debería intentar trabajar con lógica multivaluada (difusa) para obtener mejores resultados y que se necesita de un mecanismo de razonamiento para el proceso de mapeo.

Otro modelo de similitud semántica es el propuesto por Rodríguez [27][28]. Se basa en contenido de información, pero agrega ponderaciones diferentes a los atributos comunes y a los atributos diferenciadores, tanto del referente como del referenciado. Está basado en los trabajos de Tversky [29]. Lo más innovador de su trabajo es la incorporación del concepto de “intención de usuario”. Esto consiste en una implementación sencilla de contexto de uso, mediante un campo que tiene información clasificatoria. A pesar que la información de contexto es un simple atributo, los resultados mostraron que es una opción interesante a considerar.

## 2.2 Similitud de Atributos

Tversky [29] al analizar la forma en que los humanos comparan conceptos y objetos deriva que no solo los atributos iguales, sino que también los atributos diferenciadores son considerados en la evaluación. Más aún, los atributos de similitud y diferencia (*Common and Distinctive Features*) varían de acuerdo a la intención de la comparación, o con otro nombre, al contexto de significancia. Este trabajo ha influenciado a muchos otros y ha sentado bases para los modelos de similitud del futuro.

Bousquet et al. [30] mostraron un modelo usado en similitud de términos médicos. Como otras áreas del conocimiento científico con alto nivel de desarrollo de nuevos conceptos por varias personas, hay una fuerte tendencia a instanciar conceptos similares con diferentes nombres e incluso diferentes taxonomías. Este problema se repite y solo disminuye (sin desaparecer) cuando se llega a un consenso en el dominio de conocimiento, en algún futuro, y se estabilizan los cambios. Este problema es presentado por Cohen et al [31]. Constituye un problema muy importante para la construcción de ontologías y la educación de conocimiento en dominios altamente evolutivos

como ser el campo científico [32]. Otras tareas que se basan en similitud semántica sufren de los mismos problemas [33][34]. A no ser que se obtenga un modelo efectivo de integración y búsqueda de conocimiento en forma automática, va a ser difícil resolver este problema [3].

Bousquet et al. [30] sugirieron el uso de varios atributos organizados en una estructura similar a una taxonomía, y a cada diagnóstico se le asigna una combinación de conceptos. Esto es consistente con la percepción intuitiva de que problemas del mundo real tienen muchas veces causalidad compleja y que una sola taxonomía no es suficiente para analizar todas las relaciones con las que se compara un problema. El sistema dispone de una red semántica basada en las definiciones de SNOMED y se agrega una relación de causalidad. Las distancias se miden por todos los caminos entre conceptos y los arcos son pesados de acuerdo con los atributos considerados. Esto recuerda al modelo de Rips et al. [35] de ejes que describen atributos o relaciones de conceptos.

### 2.3 Basados en Teorías o en Conocimiento

Algunos modelos están basados en reglas proposicionales o de algún otro tipo que describen los objetos con los que se desea comparar. En general estas estrategias están incorporadas en sistemas mixtos donde combinan razonamiento con otros modelos. Ver [36][37][38] por ejemplos. Hay aún pocos trabajos basados en estos modelos en el área de IA, pero si hay gran cantidad de trabajos en el área de psicología cognitiva y psicolingüística que mostraron la relevancia en el pensamiento humano del razonamiento como parte del proceso de similitud [4]. Un modelo que parece promisorio en esta área es el de Razonamiento Basado en Contextos [39].

Tanto los modelos basados en teorías como los basados en atributos pueden ser utilizados como modelos de analogía para **ejemplares paradigmáticos o prototipos** ya sea que estos son dados o calculados, en estos casos la clase o categoría a la que se pertenece es representada por uno o más modelos estereotipo sobre los que se compara los atributos o las proposiciones de la teorías[40].

### 2.4 Mezcla de Taxonomía-Atributos y otros

Yamada, Inuzuka y Seki [41] crearon un modelo de similitud basado en la creencia que dados dos elementos, a mayor posibilidad de pertenecer a una clase, mayor es la similitud de ambos elementos, donde la mayor probabilidad es cuando pertenecen a la misma clase. Una métrica de similitud no existe por si misma sino que es una solución al problema en el caso considerado. Para ello se necesita un cierto punto de vista de interpretación. Los atributos comunes y diferenciadores de los objetos son medidos en relación a la cantidad de información contenida en una proposición que describe las diferencias e igualdades. Se asume que la descripción del objeto puede ser dividida en descripciones de perspectivas independientes. La similitud es entonces un promedio ponderado de las similitudes entre descripciones de los objetos de cada perspectiva. Esto recuerda mucho al modelo basado en contextos [39]. La métrica de similitud utiliza un modelo de probabilidad condicional del contenido de información. Luego se adapta este modelo a un modelo de clasificación basado en *k-Nearest Neighbor* (Instance Based Learning) y se obtiene un modelo que puede ser entrenado.

### 2.5 Dependencia por Grafos.

Melnik, Garcia-Molina y Rahm [42] definieron un método usado en un sistema para mapeo de esquemas de bases de datos, mediante mapeo de grafos. A partir de una función SQL2Graph se convierte la definición del esquema de base de datos a un grafo, y los grafos obtenidos son comparados por un método iterativo. Ellos se basan en la intuición básica de que elementos de diferentes grafos son similares cuando los elementos adyacentes son similares. Los cálculos se realizan sobre la base de los nombres de las etiquetas del grafo en forma recursiva hasta que los valores de las métricas y las ponderaciones se estabilizan. Si el sistema no converge se detiene luego de un número predeterminado de iteraciones.

El trabajo de Hasegawa et al. [43] mostró algunas experiencias en un motor de categorización semántica de documentos. Se utilizaron analizadores conceptuales que convierten al texto en diagramas conceptuales [67]. Luego se compararon los grafos supuestamente canonizados para identificar las similitudes.

## 3 Criticas

En general las actuales implementaciones de todos estos modelos incluyen alguna deficiencia en su concepción de la que deriva en un pobre desempeño o justificación teórica. Algunos de estos son descritos a continuación.

### 3.1 Modelos Fuera de Contexto o de Contexto Neutro

La mayoría de los modelos consideran la similitud semántica entre palabras o conceptos fuera de contexto. Estos modelos no consideran información situacional o de contexto, términos que son frases u otros símbolos conceptuales. Esta búsqueda de conceptos y categorías fuera de contexto o de lenguaje neutro es tan largo como nuestra historia y ha tenido muy poco éxito. Comenzó con los modelos aristotélicos de la realidad y que han prevalecido en el pensamiento humano por los últimos 2500 años. Paradójicamente, Aristóteles dice que él no es el

verdadero experto en el tema, y que de haber tenido los 50 dragmas para presenciar el curso de Pródicus sería capaz de decir más [44]. Más allá de que muchas de las interpretaciones de este texto suponen una cierta ironía de parte de Aristóteles, la realidad es que solo han llegado a nuestros días las ideas de Aristóteles a través de los escritos de Platón y poco o nada de las ideas de Pródicus.

En introspección vemos que la mente humana no maneja conceptos independientes del contexto [45]. Ya sea que el concepto se asocie a otros conceptos percibidos y/o memorizados [2] o a situaciones de uso, la psicología cognitiva [29][46][47] y la lingüística[48] están entrando en el consenso de que no es posible razonar o interpretar en forma inteligente con conceptos libres de contexto.

Aún para análisis de palabras fuera de contexto, los humanos tratamos de contextualizar los términos para aplicar el conocimiento disponible, ya sea considerando la información del agente que envió el mensaje, el medio por el que fue recibido, experiencias anteriores relacionadas con este, etc. Este proceso de contextualización es una parte central del proceso de percepción y de razonamiento de similitud [49], ya que cualquier mensaje percibido debe ser conectado al conocimiento existente del agente receptor [2] a los efectos de asociarlo al conocimiento tácito que este posee. Nótese que en este caso el contexto es solamente un caso especial de conocimiento acerca del concepto, sus restricciones y aplicabilidad.

Esto plantea la necesidad de desarrollar e incorporar mecanismos de evaluación conceptual asociados a contextos específicos, dentro de los modelos de información conceptual en los sistemas de IA.

### **3.2 Relaciones a considerar y Taxonomía.**

Muchos modelos están basados en experimentos limitados a relaciones de sub-clase y super-clase, también conocidos como relaciones *is-a* o hiponimia/hipernimia. Otros han incorporado relaciones de holonimia/meronimia o *has-part*. Pero algunos experimentos con muchas más relaciones y taxonomías más abiertas sugieren que tomar en cuenta más tipos de relaciones es mejor, aún cuando esas relaciones no estén clasificadas ni se use tipo alguno de ponderación, más que la simple cuenta de vértices [12][50]. Por otra parte hay varias experiencias que mostraron que las visiones tradicionales de llevar la representación de conceptos a clases, modelos ejemplares o prototipos va en contra de la percepción humana, por lo que no serían modelos adecuados para representación de sistemas y requisitos complejos. Ver en [4][51] algunas revisiones que comparan contra muchos experimentos en el área de psicología cognitiva.

Stuckenschmidt [72] en sus conclusiones mostró como su sistema tuvo problemas para mapear conceptos cuando no se podían subsumir. Esto mostró o que el sistema carecía del conocimiento necesario, o que muchas representaciones no tenían suficiente estructura taxonómica para representarlos, o finalmente que usar solo información taxonómica no es suficiente para representar la similitud a efectos de resolver la integración de la información. Hakimpour y Geppert [25] mostraron conclusiones similares. Otros ejemplos sugieren también las limitaciones de las taxonomías como modelo.

Guarino y Welty [52] mostraron varios problemas de implementación de ontologías y sugieren la ventaja de reutilizar algunos modelos provenientes de la escuela ontológica de la filosofía. Las ontologías que se suponen deben traer orden y estructura a ciertos dominios de conocimiento son generalmente confusas y de pobre estructura taxonómica, muchas veces debido a un uso incorrecto de las relaciones *is-a*. Algunos modelos tratan de resolver esto separando los tipos de razonamiento, muchas veces en niveles de abstracción, o en tipos de teoría, como ser sintáctico, estructural y semántico. De esta manera separan los Tbox y los Abox en sus modelos de lógica y representación [53].

Pero como la información semántica es difícil de definir y como es representada por conceptualizaciones de especificaciones más o menos concertadas entre varios individuos, pero no universales, mediante un consenso amplio, son muy difíciles de obtener. Un ejemplo simple sobre lo difícil de la universalidad de los conceptos se puede ver al comparar entre varios diccionarios de cualquier lenguaje las definiciones de varios términos, especialmente términos utilizados con significados especiales en dominios de conocimiento específicos.

Entonces las ontologías en su interpretación usual de estructuras taxonómicas arbóreas, pueden ser usadas solamente en pequeños dominios con alcance limitado a donde se pueda obtener el consenso entre los diferentes actores. Pero algunos dominios como ser el lenguaje natural, son de por sí imposibles de ser capturados por la conceptualización simple de una taxonomía de este estilo [54]. Estos problemas de obtener consenso en las definiciones de clases se extienden al de obtención de ejemplares y de prototipos ya que en cualquiera de estos casos se debe hacer un compromiso sobre los mismos.

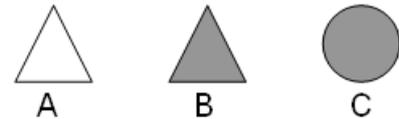
Por lo tanto estos niveles de representación (léxica, estructural y semántica) no debieran ser separados ya que están fuertemente interrelacionados. Los términos se clusterizan alrededor de conceptos, y la información sintáctica está

fuertemente relacionada a la interpretación semántica. El lenguaje natural captura mucha información contextual y situacional junto al concepto central. Al mismo tiempo dos frases en lenguaje natural referidas al mismo episodio no tendrán la misma representación ya que diferentes términos y especialmente *near-synonyms* tendrán diferente información situacional [54], haciendo imposible manejar todas las combinaciones en cualquier modelo computacional.

### 3.3 Información tácita y desigualdad triangular.

En el ejemplo de Lin [19] sobre similitud de formas geométricas, se muestra el problema de la desigualdad triangular. Considerando que la similitud es inversa a la distancia semántica, el problema entonces es el siguiente:  $dist(A,C) \leq dist(A,B) + dist(B,C)$ . En otras palabras, si A es similar a B en una cantidad a, y B es similar a C en una cantidad b, entonces A debe ser similar a C en una cantidad c tal que  $c \leq a+b$ . Y esto debería tener sentido aún para comparaciones difusas como ser {muy, mucho, un poco, etc.}. Por ejemplo dado los siguientes elementos A, B y C.

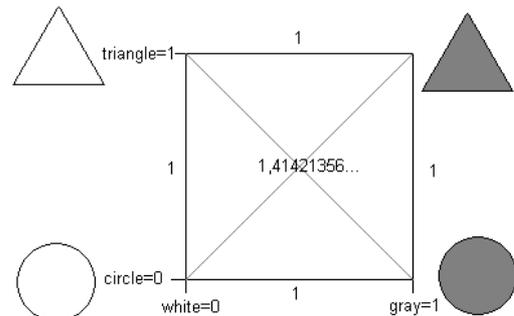
A y B tienen forma similar, pero B y C tienen similar sombreado. ¿Qué tan similares son A y C? Da la impresión que la similitud no es transitiva cuando se aplica a conjuntos disjuntos de atributos. En el ejemplo se ve que la intuición de mantener la desigualdad triangular no tiene sentido ya que lo que se quiere comparar tiene tan pocos atributos y tanta diferencia entre ellos, que el conjunto de atributos compartidos por todos los elementos es el conjunto vacío.



Para la mayor parte de los casos como este, la desigualdad triangular no tiene sentido a no ser que se identifique un conjunto de atributos relevantes para comparar conceptos dentro de un determinado contexto. Entonces el problema de la desigualdad triangular depende fuertemente del contexto de aplicación y la intención de la comparación entre conceptos, concuerda con los trabajos de Tversky en similitud humana [29].

Esta desigualdad toma sentido si se le agrega un modelo topológico o de coordenadas para comparar los atributos. Si asumimos un conjunto de coordenadas para cada atributo del concepto, donde la forma {circle, triangle} y el color {white, gray} son arreglados en ejes perpendiculares, entonces se puede derivar una medida de similitud, aún cuando la similitud intuitiva puede parecer sin sentido. Los métodos de resolución de problemas arbitrarios en humanos consisten en encontrar el contexto adecuado para la interpretación y comparación de los conceptos [48]. En el siguiente modelo, la similitud total quiere decir distancia cero, mientras que la diferencia en uno o dos atributos será de 1 o 1,41421356... (raíz de 2) respectivamente.

En este caso estamos comparando formas geométricas simples, pero los humanos como agentes cognitivos pueden derivar otros atributos y relaciones no triviales, como ser posición temporal o espacial, rugosidad de la superficie, o cualquier otra información que pueda ser derivada mediante razonamiento sobre la base de los símbolos y la situación analizada. Un agente cognitivo artificial debería poder operar en forma similar.



Una comparación por similitud es un proceso que devolverá un resultado útil si es relevante, y existe información para desambiguar el problema, de otra manera no se puede asegurar nada sobre la situación. Esto es especialmente cierto cuando los atributos similares y diferenciadores no pueden ser determinados con antelación (como en el caso de la presencia de un nuevo estímulo nunca antes percibido por el agente).

Al considerar similitud conceptual debemos incluir la información tácita ya que el token (símbolo) del concepto no explica nada acerca de su significado.

### 3.4 Contenido de Información Ponderada, Modelos Difusos y Probabilísticos

Como primer tema a notar es que la mayoría de estos modelos se han probado contra conjuntos muy reducidos de conceptos, en varios casos comparados contra el trabajo de Miller y Charles [10]. Entonces estos modelos son relevantes solo si lo que se compara es información fuera de contexto. Pero está implícito anteriormente, los problemas del mundo real deben ser conceptualizados en un contexto determinado por lo que estos modelos en general no parecen ser aplicables a problemas del mundo real.

Estos sistemas tienden a caer en algunos problemas recurrentes, como el identificado por Resnik [15] relacionado con el uso incorrecto del significado de una palabra. Significa que la polisemia debe ser resuelta antes de elegir el concepto correcto, de lo contrario la similitud se basaría en un error de categorización [1][55]. Esto mostró la importancia del contexto en el que las palabras son interpretadas [48], para asegurar que el concepto correcto se mapeado al término usado. Sin un método para extraer el significado correcto de una palabra, Resnik sugiere usar

todos los conceptos asociados con los dos términos a comparar de acuerdo a una medida de relevancia [56]. El sugiere la combinación de varios métodos de abducción de conocimiento y desambiguación de palabras a efectos de identificar el significado correcto de una palabra [57].

Por otra parte, los modelos basados en contenido de información consideran solamente el largo o la cantidad de conceptos sin considerar su significado. Por lo que, extender estos criterios de información conceptual es bastante irreal [58]. Lo dicho se cumple también para muchos trabajos basados en extensiones del de Resnik. Dekang Lin [19] mostró que la probabilidad debe ser calculada por adelantado significando un costo importante cuando hay que cambiar información del sistema.

Otros modelos de similitud basados en modelos probabilísticos y/o difusos usan esquemas taxonómicos o basados en ejemplos paradigmáticos o prototipos pero cuya similitud se calcula con modelos bayesianos y/o fuzzy. El primer tema es que no hay ninguna evidencia de que el razonamiento humano se base en modelos estadísticos, más bien todo lo contrario. No es el caso de modelos difusos donde hay muchos elementos del razonamiento humano. Esto de todas formas no sería ningún problema pues lo que estamos tratando de modelar es similitud para agentes computacionales. Pero lo cierto es que estos modelos tienen buenas respuestas para grandes cantidades de casos y no pueden responder bien para situaciones límite o para conjuntos de datos con los que no han sido entrenados. Las estadísticas de estos modelos dan buenas correlaciones cuando la cantidad de datos de entrenamiento es grande. Ver [36] como un ejemplo.

Todos estos modelos son en realidad el resultado de falta de conocimiento asociado a los conceptos que se están manejando. Esto se ve reflejado en el uso del máximo logaritmo de la probabilidad condicional de encontrar un elemento común que subsume a ambos conceptos. Este cálculo es realizado considerando toda la taxonomía incluyendo conceptos no relevantes.

### **3.5 Modelos basados en Nearest k Neighbors**

Estos modelos tienen buenas aproximaciones a los resultados pero no permiten manejar cualquier tipo de excepción, como las que ocurren bastante frecuentemente en el mundo real y en el Procesamiento de Lenguaje Natural (NLP) o en la Abducción Automática de Conocimiento (AKE). Los modelos basados en  $k$  *Nearest Neighbor* asumen un conjunto fijo y conocido de atributos para la clasificación de los elementos. Esto significa que no podemos trabajar en elementos cuya cantidad y calidad de atributos depende del contexto y/o razonamiento sobre el concepto. Ver el ejemplo de Yamada et al.[41].

### **3.6 Criterios de Similitud Humana**

Como se indicó antes, algunos autores han considerado a la evaluación humana como paradigmática y si bien no se puede reproducir, se busca que los algoritmos imiten el desempeño y los resultados de estos experimentos. Hay que diferenciar aquellos que usan los modelos para tomar ideas para sistemas computacionales con capacidad de similitud semántica, y aquellos otros que los utilizan buscando desarrollar un modelo del razonamiento y pensamiento humano.

Tenemos la experiencia de Tversky y Kahneman [59] mostrando como los humanos somos evaluadores bastante malos que utilizan tan solo unas pocas estrategias cuando se dispone de información incierta e incompleta, y el nivel de similitud expresado es relativo a la intención y a la forma en que se expresa la solicitud.

Las posiciones considerando lo opuesto son pocas y hablan de casos aislados de expertos en su tema. Uno de los ejemplos es el de la carta de Pascal enviada a Fermat [60] con relación al señor de Meré, donde se describe en esta persona un conocimiento muy preciso sobre probabilidad aprendido de la experiencia sin el conocimiento de las matemáticas.

También la evaluación humana tiene desvíos debido a sus limitados recursos cognitivos [59], por lo tanto el uso de estos modelos sin cuestionarse la aplicabilidad y que el modelo en realidad no solamente esté reproduciendo una capacidad humana sino también sus limitaciones, los hace parecer dudosos como modelos computacionales. Muchos estudios en conocimiento humano deben manejar las limitaciones fisiológicas de la mente para manejar conjuntos (chunks) de estructuras conceptuales [61][62]. Estas son limitaciones que restringen a los modelos cognitivos humanos pero que ciertamente no queremos que se propaguen a los modelos de razonamiento computarizado.

Sobre la base de los trabajos de Tversky [29] podemos ver en tareas de diferenciación o de similitud, los humanos tienden a considerar atributos relevantes a efectos de clasificar los objetos, e ignorar otros atributos no relevantes para el proceso de clasificación. Esto sugiere el uso de algún conocimiento o meta-conocimiento para manejar los problemas de similitud. La percepción de similitud puede verse afectada por la forma como se presente el problema de similitud y otros aspectos subjetivos de los cuales no se tiene la seguridad que fueran controlados en esos experimentos [63]. Sumados estos factores parece muy discutible la validez de algunos de los experimentos citados

y el uso de estos, a no ser para obtener modelos que sirven para analizar y modelizar el pensamiento humano.

Esto abre una cierta duda en varios de los modelos computacionales basados en teorías de la psicología cognitiva, sugeridos sin el suficiente análisis de validez dentro del campo de la IA. Muchos avances en la IA están relacionados a investigaciones de la psicología cognitiva, pero la aplicabilidad de estos resultados debe fundamentarse también en la aplicabilidad de las precondiciones.

Los modelos basados en contenido de información fueron algunas veces directa o indirectamente inspirados por el trabajo de Krumhansl [14]. Este modelo está relacionado a algunos de los resultados de Tversky sobre la percepción de similitud y categorización y la relevancia de los atributos de diferencia y similitud. En una taxonomía o clase muy llena, con muchos elementos para categorizar, los humanos tienden a extremar esfuerzos para diferenciar y seleccionar los elementos relevantes a efectos de acomodar los limitados recursos cognitivos humanos, y de esta manera manejar la complejidad de este conocimiento. Entonces estos modelos de densidad de información pueden ser interesantes como modelos de conocimiento humano y sus limitaciones fisiológicas, pero probablemente serán inconvenientes o irrealistas para conocimiento tratable por computadora. Por otra parte hay que notar que Corter [64] reprodujo experimentos de similitud evaluada por humanos y no pudo reproducir el efecto de Krumhansl.

#### 4 Conclusiones

La similitud semántica no debe basarse en una cantidad limitada de relaciones. Muchas veces las implementaciones de ontologías limitan estas relaciones a unos pocos tipos incluidos en la herramienta que les da soporte. Aunque recientemente varios investigadores han estado incluyendo nuevas relaciones, muchas otras son dejadas afuera. Finalmente todo tipo de relación debe ser considerado. Esto parece dar la prerrogativa a modelos basados en teorías o conocimiento donde la flexibilidad de implementar nuevas relaciones cuando las haya parece ser más sencillo. Pero también hay carencias en teorías adecuadas que determinen qué y cómo deben ser representadas estas relaciones, a no ser por algunos estudios de relaciones de mereológicas y teoría de clases.

También se necesita alguna información de explicaciones, información causal sobre los conceptos y relaciones que permita conectar la similitud por estos conceptos, ya sea por intención o por uso, u otro conocimiento causal. Varios estudios de Quillian [65], o de Collins y Loftus [66], apoyan la teoría que la mayor parte de los procesos de razonamiento en humanos puede ser simulados por redes semánticas y técnicas de *spreading activation*. De todas formas la mayor parte de los experimentos se basaron en estructuras limitadas como ontologías taxonómicas. Esto parece indicar que una solución podría estar en variaciones de modelos de grafos conceptuales, como los de Sowa [67], resolviendo algunas de sus limitaciones. Por otro parte estos modelos ofrecen explicaciones semánticas superiores a otros modelos lógicos y según su implementación permiten manejos de teorías presentenciales de explicación de la realidad [68]. Algunos de los problemas de los modelos actuales de representación de grafos conceptuales pueden verse en el trabajo de Theodorakis [69].

Euzenat [70] sugirió el uso de varios lenguajes para soportar varios modelos de manejo de conocimiento, razonamiento y otras acciones en el conocimiento compartido. Pero el problema de la similitud semántica necesita un lenguaje y un proceso de razonamiento especial para sí. Un posible modelo podría identificar la relevancia de un problema y determinar el subsistema o la interfase con la que comunicar el problema a resolver en un “servidor resolvente de problemas” específico, donde operen agentes algorítmicamente eficientes.

Euzenat, explícitamente indicó que no se debe diferenciar el conocimiento de fondo o de contexto con las anotaciones formales sobre los conceptos, pues todos forman la representación del concepto, y expresó la problemática de los lenguajes de representación con el compromiso entre expresividad y complejidad y la completitud de los demostradores de teoremas.

Una solución propuesta sería un lenguaje, extensible, lo suficientemente abierto para manejar nuevos formalismos (casi como un lenguaje de programación), con capacidad de manejo de patrones, en especial patrones situacionales. Este lenguaje debe tener previsto el pasaje de información y problemas a otros agentes más especializados.

En el manejo de Lenguaje Natural o en la abducción de conocimiento frases como “seguido”, “casi”, “en gran medida”, “la mayor parte del año”, “inmediatamente después”, etc. pueden ser encontrados en muchos términos de problemas reales [72]. Aunque imprecisos estos valores agregan valiosa información a la definición de problemas y pueden ser mapeados a rangos de valores.

La ambigüedad de estas frases no puede ser clarificada sin una fuerte investigación y el desarrollo de consenso en los modelos actuales, requiriendo grandes cantidades de esfuerzo, cuando no imposibles tareas. Además la extensión con la que deben ser evaluadas puede variar de acuerdo a la aplicación o al dominio donde son evaluadas. Pero un nuevo modelo de similitud dependiente del contexto podría considerar las diferencias entre las frases “Pablo pasó

cerca de Juan” y “Galileo pasó cerca de Júpiter”. La mayor parte de los modelos carecen de información de contexto y conocimiento del dominio de los elementos evaluados, con la excepción de Rodríguez [27], aunque en un modelo simple.

Por lo tanto, el modelo lógico de similitud debe tener bases en modelos de redes semánticas por lo que la teoría de Razonamiento Basado en Contextos al estilo de McCarthy[71][39] no sería aplicable tal como está desarrollada hoy.

Artale et al [53] sugirieron que el manejo de múltiples excepciones es un atributo fundamental para un modelo que maneje la representación y el razonamiento de similitud. Stuckenschmidt [72] sugirió que la clave para la similitud semántica es encontrar las relaciones que vinculan a las clases y objetos, y que estas relaciones tendrán diferentes modos de razonamiento y formas. Por lo tanto, no podemos limitarnos a modelos muy estructurados que limiten nuestra capacidad de razonamiento en ciertas situaciones. Entonces, el modelo debe considerar no sólo taxonomía sino también atributos, reglas y restricciones [3]. Este tipo de lenguajes tendrá la dificultad de no poder asegurar su clausura por lo que deberá estar apoyado en heurísticas y en meta-conocimiento que aseguren los procesos de razonamiento.

Cualquiera sea el modelo de representación elegido, este debe capturar información situacional, conceptual, estructural, léxica y debe estar basado en conocimiento sobre las reglas y restricciones de toda esta información. Debe, además, representar esta información en modelos de redes más complejas y flexibles que las estructuras arbóreas de muchas ontologías. Así también, debe considerar en su conocimiento no solo la información de clases, sino también las instancias dentro del modelo, como es sugerido por Rodríguez [27].

Existe la necesidad de una nueva definición de similitud y conceptos. Entonces nuestro problema de similitud semántica depende de los modelos que definamos tanto para similitud, para conceptos, dentro de una metáfora computacional que pueda manejar a todos estos junto con el conocimiento que los describe.

## 5 Agradecimientos

Se agradece los comentarios de los revisores anónimos aunque razones de espacio limitan el desarrollo de algunas de las sugerencias.

## 6 Referencias

- [1] Chi, Michelene T.H.; Creativity: Shifting Across Ontological Categories Flexibly. In T.B.Ward, S-M- Smith and J. Vaid (Eds.) *Creative Thought: An investigation of conceptual structures and processes* (pp 209-234). American Psychological Association. Washington D.C.1997
- [2] Gentner, D., & Markman, A.B. Structure mapping in analogy and similarity. *American Psychologist*, 52(1), 45-56. 1997.
- [3] Latorres, Enrique P.; Reuso de Reglas de Negocio: Una experiencia de reuso de ontologías en un dominio restringido. Master Thesis. Universidad de la República. Facultad de Ingeniería. Diciembre 2002.
- [4] Medin, D. L.. Concepts and conceptual structure. *American Psychologist*, 44(12):1469-1481, 1989.
- [5] Rota, G.C., *Indiscrete Thoughts*, F. Palombi editor, Birkhäuser, Boston-Basel-Berlin. pages 57-59, 1997.
- [6] Russel, B. *The Principles of Mathematics*, Ch XXVI, 1903.
- [7] Price, H.H. *Hume's theory of the external world*, 1940.
- [8] Osgood, C.E. The nature and measurement of meaning. *Psychological bulletin*, 49:197-237, 1952.
- [9] Rubinstein, H; Goodenough, J.B. Contextual correlates of synonymy, *Communications of the ACM*, 8(10):627-633. 1965.
- [10] Miller, G.A.; Charles, W.G. Contextual correlates of semantic similarity, *Language and Cognitive Processes*, 6(1):1-28. 1991.
- [11] Rada, R.; Mili, H.; Bicknell, E.; Blettner, M.; Development and application of a metric on semantic nets. *IEEE Transactions on Systems, Man, and Cybernetics*, 19(1):17-30, February 1989.
- [12] McHale, Michael. A comparison of WordNet and Roget's taxonomy for measuring semantic similarity. In Proc. Usage of WordNet in Natural Language Processing Systems, 1998. COLING-ACL '98 Workshop, University of Montreal, August 16, 1998.
- [13] Ramon, J.; Van Laer, W.; Bruynooghe, M.; Distance measures between atoms.
- [14] Krumhansl, C. Concerning the applicability of geometric models to similarity data: The interrelation between similarity and spatial density. *Psychological Review* 85(5):445-463. 1978.
- [15] Resnik, P; Using information content to evaluate semantic similarity in a taxonomy. In *Proceedings of IJCAI-95*, pages 448-453, Montreal, Canada. 1995.
- [16] Shannon, C.E.; A mathematical theory of communications, *Bell Systems Technical Journal*, Vol.27, pp.379-423 and 623-656, 1948.
- [17] Miller, G. WordNet: An on-line lexical database. *International Journal of lexicography*, 3(4), 1990.

- [18] Leacock, C.; Chodorow, M.; Filling in a sparse training space for word sense identification. Unpublished. 1994.
- [19] Lin, D. An Information-Theoretic Definition of Similarity, University of Manitoba, Canada. In *Proceedings of the 15<sup>th</sup> International Conference on Machine Learning*, Madison, WI.1998
- [20] Miller, G.A.; Charles, W.G. Contextual correlates of semantic similarity, *Language and Cognitive Processes*, 6(1):1-28. 1991.
- [21] Wu, Z.; Palmer, M. Verb semantics and lexical selection. In *Proceedings of the 32<sup>nd</sup> annual Meeting of the association of computational linguistics*, pages 133-138, Las cruces, New Mexico. 1994.
- [22] Levenshtein, I.V.; Binary code capable of correcting deletions, insertions, and reversals. *Cybernetics and Control Theory*, 10(8):707-710, 1966.
- [23] Jiang, J.J.; Conrath, D.W. Semantic similarity based on corpus statistics and lexical taxonomy. In *Proceedings of International Conference on research in Computational Linguistics*, Taiwan. 1997.
- [24] Budanitsky, A.; Hirst, G. Semantic Distance in Wordnet: An experimental, application-oriented evaluation of five measures. 1999
- [25] Hakimpour, F.; Geppert, A.; Resolving semantic heterogeneity in schema integration: An ontology based approach. In Chris Welty and Barry Smith, editors, *Proceedings of International conference on Formal Ontologies in Information Systems FOIS'01*. ACM Press, October 2001.
- [26] Goh, Cheng Hian; Bressan, Stephane; Madnick, Stuart; Siegel, Michael. Context interchange: New features and formalisms for the intelligent integration of information. *ACM Transaction on Information Systems*, 17(3):270-290, 1999.
- [27] Rodríguez, M.A. *Assessing semantic similarity among spatial entity classes*, PhD Thesis, University of Maine, May 2000
- [28] Rodríguez, M.A.; Egenhofer, M. J., Putting similarity assessment into context: Matching functions with user's intended operations, University of Maine.
- [29] Tversky, A. Features of similarity, *Preference, Belief and Similarity*, Selected Writings Amos Tversky, Edited by Eldar Shafir. MIT Press 2004.
- [30] Bousquet, C.; Jualet, M.C.; Chatelier, G.; Degoulet, P. Using semantic distance for the efficient coding of medical concepts. American Medical Informatics Association, AMIA 2000 Annual Symposium, Los Angeles, CA, November 2000.
- [31] Cohen, P.; CHaudri,V.; Pease, A.; Schrag, R.; Does Prior Knowledge facilitate the Development of Knowledge-Based Systems? Department of Computer Science, University of Massachusetts, 1999.
- [32] Menzies, T.; Cost Benefits of Ontologies, *Intelligence*, Fall 1999, p26-32, 1999.
- [33] Cleverdon, C.; Optimizing convenient online access to bibliographic databases, *Information Services and Use*, 4, 37-47, 1984.
- [34] Ellis, D.; Furner-Hines, J.; Willett, P. On the measurement of inter-linker consistency and retrieval effectiveness in hypertext databases, *Proceedings of SIGIR-94*, 17th ACM International Conference on Research and Development in Information Retrieval, Dublin, IE, 51-60, 1994
- [35] Rips, L.; Shoben, J.; Smith, E. Semantic Distance and the verification of Semantic relations. *Journal of Verbal Learning and Verbal Behavior*. 12:1-20. 1973.
- [36] G. Gibbon and J. Aisbett (2000) Human categorisation and its application to automatic classification. In C. Davis, T. van Gelder & R. Wales, eds., *Cognitive Science in Australia*, Causal Productions (electronic) Adelaide. 2000.
- [37] Kalish, M. and Kruschke, J. Decision boundaries in one dimensional categorization. *Journal of Experimental Psychology: Learning, Memory and Cognition* 23, 6, 1362-1377. 1997.
- [38] Kruschke, J. Base rates in category learning. *Journal of Experimental Psychology: Learning, Memory and Cognition* 22(1): 3-26. 1996.
- [39] Benerecetti, M.; Bouquet, P.; Ghidini, C.; Contextual Reasoning Distilled. *Journal of Theoretical and Experimental Artificial Intelligence*, 12(3):279-305, 2000.
- [40] Aamodt, A.; Plaza, E. Case based reasoning: foundational issues, methodological variations and system approaches., *AI Communications*. IOS Press, Vol.7:1, pp.39-59. 1994.
- [41] Yamada, Y.; Inuzuka, N.; Seki, H.; MAP classification with a similarity measure, N. Ishii (Ed.), *Proceedings of The IASTED International Conference on Artificial and Computational Intelligence (ACI 2002)*, pp. 155-160, 2002..
- [42] Melnik, S.; García-Molina, H.; Rahm, E.; Similarity Flooding: A Versatile Graph matching algorithm and its application to schema matching. In *Proceedings of 18<sup>th</sup> International Conference on Data Engineering (ICDE)*, San Jose, CA, 2002.
- [43] Hasegawa, F.; Dos Santos, E. L.; Ávila, B. C. and Kaestner C. A. A.; An Overview of Memory: Some Issues on Structures and Organization in the Legal Domain. IV Workshop on Agents and Intelligent Systems. CACIC 2003.
- [44] Plato: Cratylus (Dialog of Hemogenes and Socrates, about sayings by Cratylus)

- [45] Barsalou, L.W.; Being There Conceptually: Simulating categories in preparation for Situated Action. In N.L.Stein, P.J. Bauer, & M. Rabinowitz (Eds.) *Representation, memory and development: Essay in honor of Jean Mandler*, Mahwah, NJ:Erlbaum. October 2000.
- [46] Werner, Heinz; Kaplan, Edith; The Adquisition of word meanings: A developmental study. *Monograph of the Society for Research in Child Development*, N° 51, 1952.
- [47] Sternberg, R.J.; Powell, J.; Comprehending verbal Comprehension, *American Psychologist*, 38, 878-893. 1983
- [48] Gauker, Christopher, *Words without meaning*. MIT Press, Cambridge Massachusetts, 2003.
- [49] Sternberg; R. J.; *Cognitive Psychology*, Chapter 6 Knowledge Representation and Information Processing, p.196-220, Harcourt Brace College Publishers, 1996.
- [50] Jarmasz, M.; Szpakowicz, S. Roget's Thesaurus and Semantic Similarity. Proceedings of Conference on Recent Advances in Natural Language Processing (RANLP 2003), Borovets, Bulgaria, 212-219, September 2003.
- [51] Smith, E.E.; Medin, D.L.; *Categories and Concepts*. Cambridge, MA, Harvard University Press, 1981.
- [52] Guarino, N., and Welty, C.;Towards a methodology for ontology-based model engineering. In Proceedings of the ECOOP-2000 Workshop on Model Engineering. Available. 2000.
- [53] Artale, Alessandro; Franconi, Enrico; Guarino, Incola; Open Problems for Part-Whole Relations. 1996 international Workshop on Description Logics (DL'96), Boston MA, November 1996.
- [54] Edmonds, P; Hirst, G. Near Synonyms and Lexical Choice. *Computational Linguistics*, Volume 1, Number 1, Associations of Computational Linguistics. 2002.
- [55] Chi, M.T.H.; Roscoe, R.D. The Processes and Challenges of Conceptual Change. In Limon and Mason (Eds.) *Reconsidering Conceptual Change: Issues in Theory and Practice*, Kluwer Academic publishers, pp 3-27, 2002.
- [56] Resnik, P. Semantic classes and syntactic ambiguity. In *Proceedings of the 1993 ARPA Human Language Technology Workshop*. Morgan Kaufmann, March 1993.
- [57] Resnik, P. Semantic Similarity in a Taxonomy: An information based measure and its application to Problem of ambiguity in Natural Language. *Journal of Artificial Intelligence*, 1998.
- [58] Sowa, J. F., *Knowledge Representation: Logical, Philosophical, and Computational Foundations*, Section 2.3 Top-Level Categories. Brooks Cole Publishing Co., Pacific Grove, CA, 2000.
- [59] Tversky, A., and Kahneman, D., Judgement under uncertainty: Heuristics and biases, In *Judgement under uncertainty: Heuristics and biases*, Kahneman, Slovic and Tversky Editors. Cambridge University Press, 1982. Reprint 2001.
- [60] Pierre de Fermat. Lettre de Blaise Pascal à Pierre de Fermat. En *Oeuvres*, 29 juillet 1654.
- [61] Phillips, S., Halford, G. S., & Wilson, W. H. The processing of associations versus the processing of relations and symbols: A systematic comparison. Proceedings of the Seventeenth Annual Conference of the Cognitive Science Society, Pittsburgh,PA, 688-691, 1995.
- [62] Woods, W.; What's in a Link: Foundations for Semantic Networks, in D.G. Bobrow & A. Collins (eds.), *Representation and Understanding*, Academic Press; reprinted in, Collins & Smith (eds.), *Readings in Cognitive Science*, section 2.2. 1975.
- [63] Kahneman, D.; Tversky, A.; Choices, Values and Frames, In *Choices Values and Frames*, Daniel Kahneman and Amos Tversky Editors, Cambridge University Press, Reprinted, 2002.
- [64] Corter, J.E. Similarity, confusability, and the density hypothesis. *J. Exp. Psychol.: General*. 116:238-49. 1987.
- [65] Quillian, M.R. Semantic Memory, In Minsky, M. Ed., *Semantic Information Processing*. MIT Press, Cambridge. MA. 1968
- [66] Collins, A.; Loftus, E. A spreading activation theory of semantic processing. *Psychological Review*, 82, 407-428. 1975.
- [67] Sowa, J.F.; *Conceptual Structures: Information Processing in Mind and Machine*, Addison-Wesley, Reading, MA 1984
- [68] Horwich, Paul; *Truth*, Basil Blackwell Ltd, Oxford, 1990.
- [69] Theodorakis, M.; Analyti, A.; Constantopoulos, P.; Spyrtatos, N.; Contextualization as an Abstraction Mechanism for Conceptual Modeling. 1999.
- [70] Euzenat, Jérôme; Towards a principled approach to semantic interoperability, Proc. IJCAI workshop on Ontologies and information sharing, Seattle (WA USA), 2001.
- [71] McCarthy, J.; Generality in Artificial Intelligence. *Communications of ACM*, 30(12):1030-1035, 1987.
- [72] Stuckenschmidt, Heiner. Using OIL for Intelligent Information Integration. In V. Benjamins, A. Gomez-Perez, and N. Guarino, editors, Proceedings of the Workshop on Applications of Ontologies and Problem-solving Methods, 14th European Conference on Artificial Intelligence ECAI 2000, Berlin, Germany, Aug. 21 -- 22, 2000.