

LEC_UNIFOR no ASSIN: FlexSTS - Um *Framework* para Similaridade Semântica Textual

LEC_UNIFOR no ASSIN: FlexSTS - A Framework for Semantic Textual Similarity

Jânio Freire
Universidade de Fortaleza
janio.freire@gmail.com

Vlândia Pinheiro
Universidade de Fortaleza
vladiacelia@unifor.br

David Feitosa
Universidade de Fortaleza
davidfeitosa@gmail.com

Resumo

Desde 2012, os eventos de *Semantic Evaluation* (SemEval) propõem a tarefa de Similaridade Semântica Textual (STS) como um tema de competição, demonstrando sua relevância. Em 2016, a tarefa foi, pela primeira vez, proposta para língua portuguesa, no Workshop de Avaliação de Similaridade Semântica e Inferência Textual (ASSIN), realizado durante a conferência PROPOR 2016. Neste trabalho, apresentamos o FlexSTS - um *framework* flexível para STS que combina diversos componentes como *parsers* morfológicos e sintáticos, bases de conhecimento e lexicais, algoritmos de aprendizagem automática, e algoritmos de alinhamento e cálculo da similaridade. Para a ASSIN, FlexSTS foi instanciado em três sistemas de STS para língua portuguesa. Os resultados obtidos foram comparados com uma abordagem *baseline* que utiliza o coeficiente DICE.

Palavras chave

Similaridade Textual, Similaridade Semântica, Avaliação Semântica.

Abstract

Since 2012, Semantic Evaluation series (SemEval) propose the task of Semantic Textual Similarity (STS) as an evaluation theme, demonstrating the relevance of this research topic. In 2016, the task was first proposed to the Portuguese language, in the Workshop of Semantic Textual Similarity and Inference Evaluation (ASSIN), held during the conference PROPOR 2016. In this paper, we present the FlexSTS - a flexible framework for STS combining several components as morphological and syntactic parsers, knowledge and lexical databases, machine learning algorithms, and algorithms for alignment and similarity. For ASSIN, FlexSTS was instantiated into three STS systems for Portuguese. The results were compared with a baseline approach that uses DICE coefficient.

Keywords

Textual Similarity, Semantic Similarity, Semantic Evaluation.

1 Introdução

A tarefa de Similaridade Semântica Textual (STS) (Agirre et al, 2013) visa medir o grau de equivalência semântica entre dois textos, capturando a noção de que alguns textos são mais similares que outros. Por exemplo, o par de sentenças “A organização criminosa é formada por diversos empresários e por um deputado estadual” e “Segundo a investigação, diversos empresários e um deputado estadual integram o grupo.” devem receber um valor de similaridade mais alto que o par de sentenças “Mas esta é a primeira vez que um chefe da Igreja Católica usa a palavra em público.” e “A Alemanha reconheceu ontem pela primeira vez o genocídio armênio”. STS difere da tarefa de Inferência textual (RTE), principalmente por assumir uma equivalência bidirecional, e difere da tarefa de RTE e Paráfrase por definir uma noção de grau de similaridade, ao invés de uma decisão binária (sim/não).

Computar a similaridade textual é útil para um número crescente de tarefas de Processamento de Linguagem Natural (PLN) e Inteligência Artificial (IA), tais como sumarização (Lin e Hovy, 2003) e reuso de experiência (Albuquerque et al, 2012).

Desde 2012, os eventos de *Semantic Evaluation* (SemEval)¹ propõem esta tarefa como um tema de competição, demonstrando a relevância da mesma e um tema de pesquisa ainda em aberto. Em 2016, a tarefa foi novamente proposta para língua inglesa na edição do SemEval 2016² e, de forma inédita para língua portuguesa, no Workshop de Avaliação de Similaridade Semântica e Inferência Textual (ASSIN), realizado durante a conferência PROPOR 2016³. Tradicionalmente, a tarefa consiste em computar o grau de similaridade semântica entre duas sentenças, usando a seguinte escala: (1) Sentenças completamente diferentes, em assuntos diferentes; (2) Sentenças não relacionadas, mas que

¹ <https://en.wikipedia.org/wiki/SemEval>

² <http://alt.qcri.org/semeval2016/task1/>

³ <http://propor2016.di.fc.ul.pt>

compactam do mesmo assunto; (3) Sentenças de certa forma relacionadas, que podem descrever fatos diferentes mas compartilham alguns detalhes; (4) Sentenças fortemente relacionadas, que divergem apenas em alguns detalhes; (5) Sentenças significam exatamente a mesma coisa.

Neste trabalho, apresentamos o FlexSTS - um *framework* genérico que facilita e flexibiliza o desenvolvimento de sistemas de STS, pois combina diversos componentes como *parsers* morfológicos e sintáticos (NLP toolkits), bases de conhecimento e lexicais, algoritmos de aprendizagem automática, e algoritmos de alinhamento e cálculo da similaridade. Especificamente para avaliação no Workshop ASSIN, FlexSTS foi instanciado para língua portuguesa em três configurações (sistemas) usando o parser *Freeling* (Padró e Stanilovsky, 2012), o modelo de similaridade entre palavras HAL (Hyperspace Analog to Language) (Burgess, Livesay e Lund, 1998), a base de conhecimento Wordnet (Miller, 1995), o algoritmo de aprendizagem automática proposto em (Pedregosa et al., 2011), e o modelo de alinhamento entre termos proposto em (Han et al., 2013). Foram enviadas as execuções dos três sistemas de STS e os resultados obtidos foram comparados com uma abordagem *baseline* que utiliza o coeficiente DICE (Rohlf, 1992) de similaridade sintática entre textos. A análise de casos em que nosso melhor sistema não obteve nível de acerto desejado indiciam melhorias para trabalhos futuros.

2 FlexSTS - *Framework* para Similaridade Semântica Textual

Nesta seção apresentamos a proposta do *framework* FlexSTS, o qual define diversos componentes a serem plugados no desenvolvimento de sistemas de STS, agregando modelos e medidas de similaridade, *toolkits* e algoritmos do estado da arte, em cada etapa do processo de STS. A Figura 1 apresenta o fluxo geral do processo de STS e os diversos componentes ou plug-ins necessários.

2.1 Análise Morfológica e Sintática

Nesta etapa, dados dois textos de entrada t_1 e t_2 , é realizada a detecção das sentenças, a análise morfológica (*tokenização*, lematização, *POS Tagger*) e a análise sintática (*dependence parsing*) de ambos os textos. Inúmeros *toolkits* disponíveis podem realizar esta tarefa para diversas línguas. Em destaque, tem-se o Stanford NLP Toolkit (Toutanova et al, 2000), Open NLP (Baldrige, 2005), Freeling (Padró e Stanilovsky, 2012).

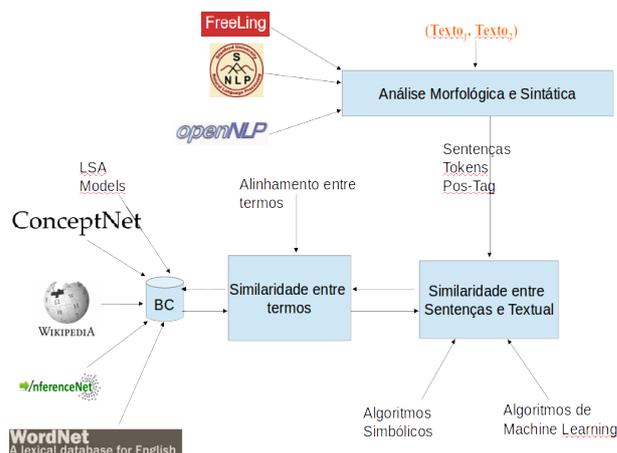


Figura 1: Fluxo do *framework*.

2.1 Análise Morfológica e Sintática

Nesta etapa, dados dois textos de entrada t_1 e t_2 , é realizada a detecção das sentenças, a análise morfológica (*tokenização*, lematização, *POS Tagger*) e a análise sintática (*dependence parsing*) de ambos os textos. Inúmeros *toolkits* disponíveis podem realizar esta tarefa para diversas línguas. Em destaque, tem-se o Stanford NLP Toolkit (Toutanova et al, 2000), Open NLP (Baldrige, 2005), Freeling (Padró e Stanilovsky, 2012).

O objetivo desta etapa é gerar, para cada texto de entrada, o conjunto de *tokens* relevantes T_{ij} de cada sentença s_{ij} . O algoritmo para a construção do conjunto T_{ij} , segue os passos listados abaixo:

1. Análise morfológica e sintática do texto;
2. Reconhecimento de palavras compostas, nomes próprios, valores numéricos, datas e expressões de tempo;
3. Aplicação de heurísticas, seguindo o trabalho de (Han et al., 2013) :
 - i. Remoção de pontuação;
 - ii. Expressões numéricas escritas por extenso são convertidas para números;
 - iii. Remoção de *stop words*.
 - iv. Referências para tempo são convertidas para o formato militar;
4. Cada *token* das classes abertas de palavras (substantivo, verbo, advérbio e adjetivo), incluindo nomes de entidades reconhecidas, como nomes próprios e abreviações, passam por um processo de desambiguação conforme definido em (Furtado et al, 2012). Nesse passo, cada termo é associado a um conceito de uma base de conhecimento.
5. Finalmente, o conjunto T_{ij} é formado pelos *tokens* e seus atributos morfológicos, lexicais, sintáticos e semânticos.

2.2 Similaridade Semântica de Palavras

A segunda etapa do processo prevê a aplicação de modelos e medidas para cálculo da similaridade entre palavras $\theta(c, c')$ e de um algoritmo para alinhamento dos termos c e c' de cada sentença s_{1i} e s_{2j} dos textos t_1 e t_2 (textos de entrada).

2.2.1 Modelos de Similaridade Semântica entre Palavras (Word Similarity Models)

O framework define a função $\theta(c, c')$ como uma função parametrizável para vários modelos e medidas de similaridade entre palavras, possibilitando agregar conhecimento adicional expresso em uma ou mais bases de conhecimento e dicionários externos, tais como Wikipedia (Witten e Milne, 2008), WordNet (Miller, 1995), ConceptNet (Liu e Singh, 2004), InferenceNet (Pinheiro et al, 2010).

Dentre os modelos do estado da arte, tem-se a LSA (*Latent Semantic Analysis*) que segue a hipótese da semântica distribucional, segundo a qual “palavras que ocorrem em contextos similares tendem a ter significados similares” (Harris, 1968). Diversas técnicas de LSA podem ser aplicadas. HAL (*Hyperspace Analog to Language*) (Burgess, Livesay e Lund, 1998) é uma variação da técnica de LSA que pode ser aplicada em matriz de coocorrência termo-termo. *Singular Value Decomposition* (SVD) (Landauer e Dumais, 1997) tem sido efetiva para melhorar medidas de similaridade entre palavras, visto que podemos selecionar os k -maiores valores singulares e reduzir para tamanho k o vetor que representa uma palavra. Por fim, a similaridade entre duas palavras é calculada pela similaridade do cosseno entre os vetores de cada palavra. Em (Han et al., 2013), tem-se uma descrição detalhada do uso do modelo HAL com SVD para língua inglesa.

O modelo de similaridade semântica inferencialista, proposto em (Pinheiro, Furtado e Albuquerque, 2014) e (Pinheiro, Pequeno e Furtado, 2010) define a *Word Inferential Similarity Measure* a qual calcula a similaridade entre dois conceitos pela interseção entre o conjunto das pré-condições [ou pós-condições] de uso dos dois conceitos, aludindo a ideia de que quanto mais as circunstâncias [ou consequências] de uso de ambos os conceitos são similares, mas as inferências em que os mesmos podem participar são similares.

Han et al. (2013) propõem uma medida de similaridade entre palavras que agrega valor da base WordNet à medida HAL.

2.2.2 Estratégias de Alinhamento entre termos

A estratégia de alinhamento é necessária para definir quais termos de cada sentença serão comparados em termos de similaridade semântica. Considere os textos de entrada t_1 e t_2 com as seguintes sentenças $\{s_{11}, s_{12}, s_{13}\}$ e $\{s_{21}, s_{22}, s_{23}\}$, respectivamente. Na etapa anterior, os conjuntos T_{11} e T_{21} com os termos das sentenças s_{11} e s_{21} foram gerados. Propõe-se então uma função de alinhamento $t_align(c)$ (Formula (1) que busca alinhar o termo c em T_{11} com um ou mais termos c' em T_{21} , de acordo com uma das seguintes estratégias:

1. *tokens* de mesma classe gramatical (*POS tag*) (p.ex. substantivo com substantivo, verbo com verbo, etc);
2. *tokens* com a mesma função sintática (p.ex. sujeito com sujeito, verbo principal com verbo principal, objeto direto com objeto direto, etc);
3. *tokens* com maior valor de similaridade semântica entre palavras;
4. todos os *tokens* com todos;

Seguindo Han et al. (2013), a estratégia 3 alinha o termo c em T_{ij} com o termo c' em T_{lj} , que tiver maior valor de similaridade semântica $\theta(c, c')$ (Formula (1)).

$$t_align_3(c) = \operatorname{argmax}_{c' \in T_{lj}} \theta(c, c'). \quad (1)$$

A flexibilidade de adotar uma dentre várias estratégias de alinhamento permite adaptar o sistema STS a um domínio ou aplicação. No entanto, argumentamos que a estratégia 1 (que utiliza o critério de *POS tag*) e a estratégia 2 (que utiliza o critério de função sintática) são mais intuitivas e linguisticamente fundamentadas, embora mais complexas.

2.3 Similaridade Semântica Textual

Na última etapa do processo, o framework prevê duas abordagens para cálculo da STS –algoritmos de aprendizagem automática e/ou algoritmos simbólicos.

A abordagem por aprendizagem de máquina preconiza o uso de algoritmos supervisionados, tais como definidos em (Chang e Lin, 2011), (Hall et al, 2009) e (Pedregosa et al, 2011), com uso de características (*features*) sintáticas, lexicais e semânticas.

Na abordagem simbólica, a intuição básica de uma medida de similaridade semântica entre textos é que, quanto mais as sentenças dos textos são similares, mais os textos são similares. Da mesma forma, quanto mais os conceitos articulados nas sentenças são similares, mais similares as sentenças

também serão. Neste sentido, a medida *SIMt* (Formula (4)) define a similaridade entre dois textos de entrada t_1 e t_2 pela média da similaridade entre as sentenças s e s' que são mais similares. Ou seja, cada sentença s de t_1 , é alinhada com a sentença s' de t_2 que lhe é mais similar.

A Formula (2) apresenta nossa função de alinhamento de sentenças $s_align(s)$, a qual, para a sentença s de t_1 (ou t_2), retorna sua contraparte s' em t_2 (ou t_1), com maior valor da medida de similaridade entre sentenças *SIMs* (Formula (3)).

$$s_align(s) = \underset{s' \in t_2}{\operatorname{argmax}} SIMs(s, s'). \quad (2)$$

A Formula (3) define a medida de similaridade entre sentenças *SIMs* entre duas sentenças s_1 e s_2 pela média ponderada do somatório das similaridades entre seus termos alinhados.

$$SIMs(s_1, s_2) = \frac{\sum_{i=1}^n \left(\sum_{j=1}^{q_i} (\theta(c, c') * P_i) \right)}{\sum_{i=1}^n (q_i * P_i)} \quad (3)$$

Onde,

- $\theta(c, c')$ é o valor da similaridade entre os *tokens* das sentenças s_1 e s_2 , de acordo com o modelo de similaridade entre palavras definido na etapa anterior;
- n é a quantidade de “tipos gramaticais” definidos na estratégia de alinhamento. Por exemplo, usando o critério de alinhamento por função sintática (estratégia 2), pode-se ter $n=3$, conforme os seguintes tipos: SUJEITO, VERBO PRINCIPAL e OBJETO;
- q_i é a quantidade de elementos em cada “tipo gramatical” i ;
- p_i é o peso do “tipo gramatical” i , permitindo, por exemplo, que a similaridade entre verbos tenha um peso maior que a similaridade entre objetos diretos.

Finalmente, a Formula (4) calcula a similaridade semântica entre dois textos de entrada t_1 e t_2 , com p e k sentenças, respectivamente.

$$SIMt(t_1, t_2) = \frac{\sum_{s \in t_1} (SIMs(s, s_align(s)))}{2 * |t_1|} + \frac{\sum_{s \in t_2} (SIMs(s, s_align(s)))}{2 * |t_2|} \quad (4)$$

Pinheiro, Furtado e Albuquerque (2014) apresentam um exemplo ilustrativo de uso das fórmulas acima.

3 Sistemas STS para ASSIN

O framework FlexSTS foi usado para instanciar três sistemas para STS na língua portuguesa, cujos resultados foram submetidos à avaliação no Workshop de Avaliação de Similaridade Semântica e Inferência Textual (ASSIN), realizado durante a conferência PROPOR 2016. A seguir serão explanadas a configuração de cada sistema e do sistema *baseline*. Ao final, os resultados e uma discussão dos mesmos serão apresentados.

3.1 STS_MachineLearning

O sistema *STS_MachineLearning* aplicou uma abordagem híbrida para cálculo da STS - aprendizagem automática usando dois atributos (*features*) - similaridade entre palavras pelo coeficiente DICE e similaridade entre palavras pela WordNet. A configuração do sistema está descrita na Tabela 1.

Etapa	Componente/Modulo	Ferramenta
Análise Morfológica/Sintática	POS <i>Tagger</i> Lematização	FreeLing
Similaridade Semântica de Palavras	Coeficiente DICE	Ver 3.1.1.1
	WordNet	Ver 3.1.1.2
Similaridade Semântica textual	Aprendizagem Automática	Ridge Regression Model

Tabela 1: Configuração do sistema STS_MachineLearning.

3.1.1 Modelo de Aprendizagem de Máquina

No cálculo de STS foi usado o algoritmo *ridge regression model* (Pedregosa et al, 2011), um modelo de regressão com $\alpha = 1.0$ e um resolvidor automático que seleciona o peso de uma coleção dependendo do tipo de dado. Esses algoritmos foram usados em (Sultan et al, 2015), campeão da tarefa de STS no SemEval 2015. O treinamento do algoritmo *ridge regression model* foi realizado com o *dataset* de treinamento disponibilizado na ASSIN. A seguir detalhamos os cálculos das duas *features* usadas para caracterizar o conjunto de exemplos.

3.1.1.1 Feature DICE

Esta *feature* representa a similaridade semântica textual entre os dois textos (exemplos) calculada pela Fórmula (4) usando uma função DICE, a qual considera o coeficiente de Dice tradicional (Rohlf, 1992), caso de números e de pronomes correspondentes. Neste caso $\theta(c,c') = \text{DICE}(c, c')$, definida na Formula (5).

$$\text{DICE}(c, c') = \begin{cases} 1, & \left(\begin{array}{l} \text{isNumber}(c) \wedge \text{isNumber}(c') \wedge c=c' \\ \text{isCorrespondingPronouns}(c, c') \\ \text{diceCoeficiente}(c, c') > 2/3 \end{array} \right) \\ 0, & \text{do contrario} \end{cases} \quad (5)$$

Onde,

- *isNumber(c)* retorna verdadeiro se o termo *c* é um número;
- *isCorrespondingPronouns(c,c')* verifica se os termos *c* e *c'* são pronomes correspondentes. Por exemplo, para os pronomes “eu” e “me” retorna verdadeiro;
- *diceCoeficiente(c,c')* calcula o coeficiente de Dice entre os termos *c* e *c'*, conforme definido em (Rohlf, 1992).

3.1.1.2 Feature WNET

Esta *feature* representa a similaridade semântica textual entre os dois textos (exemplo) calculada pela Fórmula (4) usando conhecimento da WordNet para calcular a similaridade entre palavras, conforme Formula (6):

$$\text{WNET}'(c, c') = 0.5 e^{\alpha D(c, c')} \quad (6)$$

Onde,

- $D(c, c')$ é uma função de distância entre os termos na base WordNet (distancia ≤ 4 links), calculado conforme segue:
 - 0, caso os termos pertençam ao mesmo conjunto de sinônimos (*synset*);
 - 1, nos seguintes casos:
 - uma palavra é hiperonímia direta da outra.
 - um adjetivo tem uma relação direta do tipo *similar to* com outro.
 - uma palavra é uma forma derivacional da outra.
 - 2, nos seguintes casos:

- uma palavra é 2 links de hiperonímia indireta da outra.
 - um adjetivo é 2 links *similar to* com outro.
 - uma palavra é cabeça (head) do glossário da outra, ou sua hiperônima direta, ou uma das suas hipônimas diretas.
- α , parâmetro de normalização definido em (Han et al., 2013) e fixado em 0,25.

A versão utilizada da WordNet foi a versão 3.0 em inglês e foi realizada a tradução dos corpus da ASSIN (português-inglês) pelo Google Tradutor.

3.2 STS_HAL

O sistema STS_HAL aplicou somente a abordagem simbólica para cálculo da STS, usando o modelo HAL de similaridade entre palavras e a estratégia de alinhamento por termos com maior similaridade (estratégia 3). A configuração do sistema STS_HAL está descrita na Tabela 2.

Etapa	Componente/Modulo	Ferramenta
Análise Morfológica/Sintática	POS Tagger Lematização	FreeLing
Similaridade Semântica de Palavras	Modelo (HAL+SVD)	Ver 3.2.1
	Estratégia de alinhamento	<i>t_align</i> ₃ (Formula (1))
Similaridade Semântica textual	Algoritmo matemático STS	Formulas (2), (3) e (4)

Tabela 2: Configuração do sistema STS_HAL.

3.2.1 Modelo de Similaridade HAL

Foi usada a variação da técnica LSA chamada HAL (*Hyperspace Analog to Language*) (Burgess, Livesay e Lund, 1998) que constrói a matriz de coocorrência termo-termo. Para a construção da matriz, foi usado o corpus CETENFolha⁴ - um corpus de cerca de 24 milhões de palavras em português brasileiro, com base nos textos do jornal Folha de S. Paulo que fazem parte do corpus do Núcleo Interinstitucional de Linguística Computacional (NILC), da USP/São Carlos.

Por questões de performance, foram selecionados os 24000 termos que mais ocorrem no corpus, das classes abertas de palavras (substantivos, verbos, adjetivos e advérbios). Neste

⁴ <http://www.linguateca.pt/cetenfolha/>

vocabulário não existem nomes próprios. A frequência de coocorrência entre os 24000 termos foi contada em uma janela de tamanho fixo que passa por todo o corpus. O tamanho de janela utilizado foi ± 4 , pois foi o que obteve melhor resultado em (Han et al., 2013). Por fim, foi aplicado a estratégia de SVD (*Single Value Decomposition*) de (Baglama e Reichel, 2015), e selecionados os $k=300$ maiores valores singulares. Assim, o tamanho do vetor que representa as palavras foi reduzido de 24000 para 300. A similaridade entre os termos foi calculada utilizando a função cosseno entre os vetores.

3.3 STS_WORDNET_HAL

O sistema STS_WORDNET_HAL aplicou somente a abordagem simbólica para cálculo da STS, o modelo HAL de similaridade entre palavras e a estratégia de alinhamento por termos com maior similaridade (estratégia 3). Como conhecimento adicional, adicionou informação da WordNet no cálculo da similaridade, a exemplo do trabalho de (Han et al., 2013). A configuração do sistema STS_WORDNET_HAL está descrita na Tabela 3.

Etapa	Componente/Modulo	Ferramenta
Análise Morfológica/Sintática	POS <i>Tagger</i> Lematização	FreeLing
Similaridade Semântica de Palavras	Modelo (HAL+SVD)	Ver 3.2.1
	Estratégia de alinhamento	t_align_3 (Formula (1))
	Base de Conhecimento - WordNet	Ver 3.3.1
Similaridade Semântica textual	Algoritmo matemático STS	Formulas (2), (3) e (4)

Tabela 3: Configuração do sistema STS_WORDNET_HAL.

3.3.1 HAL + Conhecimento da WordNet

À medida de similaridade entre palavras $\theta(c, c')$ foi adicionado conhecimento da base WordNet, seguindo (Hal et al., 2013). As Fórmulas (7a) e (7b) apresentam este cálculo, onde:

- $\theta(c, c') = HAL(c, c')$ (ver 3.2.1);
- *usaDice* é um parâmetro que indica se, em caso valor $\theta(c, c')$ nulo ou zerado, deva-se usar o valor da função DICE;
- $DICE(c, c')$, conforme definido em Formula (5);

- $WNET'(c, c')$, conforme definido em Formula (6).

$$WNET(c, c') = BASIC(c, c') + WNET'(c, c') \quad (7a)$$

$$BASIC(c, c') = \begin{cases} \theta(c, c') & \theta \neq \text{nulo} \\ DICE(c, c') & \text{usaDice} = \text{verdadeiro} \wedge (\theta = \text{nulo} \vee \theta(c, c') = 0) \\ 0 & \text{do contrario} \end{cases} \quad (7b)$$

3.4 STS_Baseline

O sistema *STS_Baseline* foi usado neste trabalho apenas como referência inicial de avaliação, visto que, antes da ASSIN, inexistia estado da arte para STS em língua portuguesa. Nossa proposta foi utilizar o coeficiente de similaridade DICE (conforme definido em 3.1.1.1), como sistema *baseline* para a tarefa de STS.

3.5 Resultados e Discussão

A tabela 4 apresenta os resultados da medida de correlação de *Pearson* dos três sistemas STS (runs), enviados para ASSIN, após execução no *dataset* de teste para Português-Brasileiro (PT-BR) e Português-Portugal (PT-PT). Nosso melhor sistema foi o STS-MachineLearning em ambos os *datasets* PT-BR e PT-PT. Na última linha da Tabela 4, apresentamos os resultados do sistema *baseline*, que obteve melhor performance que qualquer um dos sistemas avaliados para PT-PT.

Sistema	PT-BR	PT-PT
STS_MachineLearning	0,62	0,64
STS_HAL	0,56	0,59
STS_WNET_HAL	0,61	0,63
STS_Baseline	0,60	0,69

Tabela 4: Resultados dos sistemas STS desenvolvidos a partir do *framework* FlexSTS.

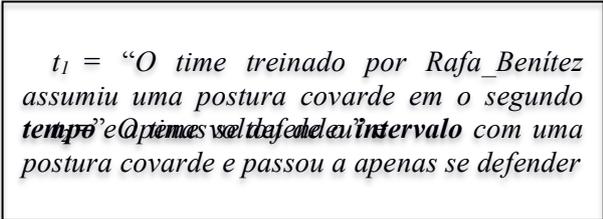
A seguir elencamos duas dificuldades importantes enfrentadas na construção dos sistemas de STS submetidos à ASSIN:

- No sistema STS_HAL, a matriz de coocorrência termo-termo gerada era muito esparsa, implicando em pouca relevância do cálculo da similaridade pela HAL. Atribuiu-se como causa o tamanho do corpus e tamanho dos textos do corpus;
- O uso da versão em Inglês da WordNet com a necessidade de solução de tradução português-inglês dos corpus ASSIN pode

ter prejudicado o desempenho dos sistemas que utilizam esta base.

O uso do sistema *baseline* pelo coeficiente DICE permitiu constatar que uma medida simples de similaridade sintática obteve resultado significativo em relação aos corpus PT-BR (0,60) e PT-PT (0,69). Em apenas 211 casos do corpus Gold Standard ASSIN, o valor absoluto da diferença entre o valor da similaridade DICE e o valor GOLD foi superior a 2 ($|DICE-GOLD| > 2$). No demais casos (1935), estes valores são muito próximos. Portanto, conclui-se que os corpus ASSIN possuem uma similaridade lexical alta, dificultando a influência de conhecimento semântico à tarefa de STS.

Analisando alguns casos em que o sistema STS_MachineLearning obteve melhor resultado comparado com a solução *baseline* (DICE), identificamos que conhecimento semântico agregou valor à tarefa. Por exemplo, para o par de texto t_1 e t_2 na Figura 2, o sistema STS_MachineLearning apresentou valor de similaridade mais correlato ao valor GOLD, pois encontrou valor de similaridade entre as palavras “*intervalo*” e “*tempo*”.



$t_1 =$ “O time treinado por Rafa Benítez assumiu uma postura covarde em o segundo **tempo** e ~~optima~~ ~~se~~ ~~na~~ ~~tarefa~~ ~~de~~ ~~o~~ ~~intervalo~~ com uma postura covarde e passou a apenas se defender

Figura 2: Exemplo de textos com uso de conhecimento da WordNet.

4 Trabalhos Relacionados

Destacam-se, como estado da arte, os sistemas campeões da tarefa de STS das edições do SemEval 2013, 2014, 2015.

No SemEval 2013, o sistema campeão foi o submetido pela equipe denominada UMBC (Han et al., 2013). Esse sistema consiste de uma abordagem que agrega conhecimento semântico de uma matriz HAL e da WordNet, além de aplicar uma estratégia de alinhamento e penalização, que determina um conjunto de critérios para um mal alinhamento, e valores e a serem descontados para cada tipo de mal alinhamento. O resultado médio da correlação de Pearson foi 0.6181, para língua inglesa.

Em 2014, a equipe vencedora foi a ECNU (Zhao et al., 2014) que utilizou uma abordagem de aprendizagem de máquina com vários algoritmo e 72 *features*. O algoritmo que obteve melhor resultado foi o *Gradient Boosting*. O resultado

médio da correlação de Pearson foi 0,8414, também para língua inglesa.

O sistema campeão da edição de 2015 foi o apresentado em (Sultan et al, 2015) que propôs uma abordagem de aprendizagem de máquina utilizando o algoritmo *Ridge Regression Model*. As características (*features*) definidas para representar o problema baseiam-se na similaridade entre as sentenças, calculada por uma função que usa uma representação vetorial, criada a partir da matriz HAL, de uma base de paráfrase (Ganitkevitch et al., 2013) e da árvore de dependência sintática. Este sistema obteve performance de 0,8015 (correlação de Pearson).

5 Conclusão

Neste trabalho apresentamos a proposta do *framework* FlexSTS, o qual define diversos componentes a serem selecionados para o desenvolvimento de sistemas de STS, agregando modelos e medidas de similaridade, *toolkits* e algoritmos do estado da arte, em cada etapa do processo de STS.

FlexSTS foi instanciado em três sistemas: (1) **STS_MachineLearning** - abordagem híbrida para cálculo da STS com aprendizagem automática usando dois atributos (*features*) - similaridade entre palavras pelo coeficiente DICE e similaridade entre palavras pela WordNet; (2) **STS_HAL** - abordagem simbólica que usa basicamente o modelo de similaridade de palavras da *Latent Semantic Analysis* (LSA); (3) **STS_WORDNET_HAL** - uma abordagem também simbólica que agrega conhecimento da WordNet à similaridade pela LSA. Os sistemas foram testados nos *datasets* de teste disponíveis na ASSIN para Português brasileiro (PT-BR) e Português de Portugal (PT-PT). Nosso melhor sistema foi o STS-MachineLearning com resultado para o PT-PT de 0,64 (correlação de Pearson). Os principais problemas foram a esparsidade da matriz de coocorrência termo-termo construída a partir do corpus CETEMFolha e o uso da WordNet em inglês. Um resultado importante foi o desempenho do sistema *baseline* pelo coeficiente de DICE, eu obteve 0,69 para o *corpus* PT-PT, indiciando que os corpus possuem alta similaridade lexical.

A análise dos resultados, dos problemas enfrentados e de erros do sistema indicam os seguintes trabalhos futuros: criação de mais cenários de testes com diversificação de algoritmos de machine-learning e novas *features*; construção de nova matriz HAL a partir de um corpus mais robusto na língua portuguesa; agregação de conhecimento da Wikipedia e InferenceNet.

Referências

- Agirre, E, et al. 2013. Sem 2013 shared task: Semantic textual similarity. *SEM 2013: The Second Joint Conference on Lexical and Computational Semantics.
- Albuquerque, A., Pinheiro, V., Leite, T. 2012. Reuse of Experiences Applied to Requirements Engineering: An Approach Based on Natural Language Processing In: Proceedings of the 24th International Conference on Software Engineering and Knowledge Engineering, SEKE 2012, São Francisco, CA.
- Baglama, J. e Reichel, L. (2015). irlba: Fast Truncated SVD, PCA and Symmetric Eigen decomposition for Large Dense and Sparse Matrices. R package version 2.0.0
- Baldrige, Jason. "The opennlp project." URL: <http://opennlp.apache.org/index.html>,(accessed 2 February 2012) (2005).
- Burgess, C., K. Livesay, e K. Lund. 1998. Explorations in context space: Words, sentences.25:211 – 257.
- Chang, Chih-Chung e Lin, Chih-Jen. 2011. LIBSVM: A library for support vector machines. ACM Transactions on Intelligent Systems and Technology, 2:27:1-27:27.
- Furtado, V., Pinheiro, V., Freire, L., Ferreira, C. (2012) Knowledge-Intensive Word Disambiguation via Common-Sense and Wikipedia In: 21st Brazilian Symposium on Artificial Intelligence. SBIA 2012, 2012, Curitiba, PR. Lecture Notes in Artificial Intelligence (LNAI).
- Ganitkevitch, Juri, BenjForcada,amin Van Durme, and Chris Callison-Burch. 2013. PPDB: The Paraphrase Database. In Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics, NAACL-HLT '13, pages 758-764, Atlanta, Georgia, USA.
- Hall, Mark, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, and Ian H. Witten. 2009. The WEKA Data Mining Software: An Update. SIGKDD Explorations, 11(1).
- Han, L., Kashyap, A.L., Finin, T., Mayfield, J., and Weese, J. 2013. UMBC_EBIQUITY-CORE: Semantic Textual Similarity Systems. In Proceedings of the Second Joint Conference on Lexical and Computational Semantics.
- Harris, Z. Mathematical Structures of Language. Wiley, New York, USA. (1968).
- Landauer, T. K., & Dumais, S. T. (1997). A solution to Plato's problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. Psychological review, 104(2), 211.
- Lin, Chin-Yew e Hovy, Eduard. 2003. Automatic evaluation of summaries using n-gram cooccurrence statistics. In Proceedings of Human Language Technology Conference (HLT-NAACL 2003),Edmonton, Canada,
- Liu, H., Singh, P. 2004. ConceptNet: A Practical Commonsense Reasoning Toolkit. BT Technology Journal, Volume 22(4). Kluwer Academic Publishers.
- Miller, George A. 1995. Wordnet: a lexical database for english. Communications of the ACM, 38(11):41.
- Padró, Lluís e Evgeny Stanilovsky. 2012. Freeling 3.0: Towards wider multilinguality.
- Pedregosa, Fabian et al., 2011. Scikit-learn: MachineLearning in Python. Journal of Machine Learning Research, vol. 12, pages 2825-2830.
- Pinheiro, V., Furtado, V., Albuquerque, A. 2014. Semantic Textual Similarity of Portuguese-Language Texts: An Approach Based on the Semantic Inferentialism Model. In: Baptista, J et al. (eds.): PROPOR 2014, Lecture Notes in Computer Science Volume 8775, 2014, pp 183-188 Springer, Heidelberg
- Pinheiro, V., Pequeno, T., Furtado, V. 2010. Um Analisador Semântico Inferencialista de Sentenças em Linguagem Natural. Linguamática. ISSN: 1647-0818. Vol.2. Num.1, pp. 111-130.
- Pinheiro, V., Pequeno, T., Furtado, V., Franco, W. InferenceNet.Br: Expression of Inferentialist Semantic Content of the Portuguese Language. In: T.A.S. Pardo et al. (eds.): PROPOR 2010, LNAI 6001, pp. 90-99. Springer, Heidelberg.
- Rohlf, F. J. Numerical taxonomy and multivariate analysis system. Version 1.70. New York: [s.n.], 1992. 470 p.
- Sultan, Md Arafat, Steven Bethard, Tamara Sumner. 2015. Dls@cu: Sentence similarity from word alignment and semantic vector composition. pp. 148 – 153. In SemEval 2015.
- Toutanova, Kristina ,Klein ,Dan , Manning, Christopher,Morgan, William ,Rafferty, Anna , and Galley, Michel . . Stanford log-linear part-of-speech tagger. <http://nlp.stanford.edu/software/tagger.shtml>.(2000)
- Witten, I., Milne, D. 2008. An effective, low-cost measure of semantic relatedness obtained from Wikipedia links. Proceeding of AAAI Workshop on Wikipedia and Artificial Intelligence: an Evolving Synergy (pp. 25-30). Chicago: AAAI Press.

Zhao, Jiang, Tian Tian Zhu, Man Lan. 2014.
ECNU: One Stone Two Birds: Ensemble of
Heterogenous Measures for Semantic
Relatedness and Textual Entailment. In

Proceedings of the 8th International Workshop
on Semantic Evaluation (SemEval-2014).