

# Análise dos Algoritmos base do Graphia

André Felipe  
Karen Teixeira

# Visão Geral

O projeto tem como objetivo analisar os resultados de extração de relações e entidades usando os algoritmos de pos-tagger e análise sintática baseados no Graphia, para atribuir semântica a essas relações.

# Texto Escolhidos

- Rio de Janeiro
- Rede Globo
- Roberto Marinho

Inspiração: Escolher textos em português que envolvessem uma pessoa, um lugar e uma empresa que possuísem relações. No nosso trabalho, utilizamos textos retirados da Wikipedia em português, para que pudéssemos extrair automaticamente os textos.

# Passo a Passo da extração das árvores de análise sintática

## Passo 1:

Recuperação de um dump da Wikipedia no dia 04/11/2014 e usamos o script em python chamado WikiExtractor.py criado por um italiano que pega o dump que vem em XML e transforma num texto limpo, descartando qualquer tipo de informação ou anotação presentes na página da Wikipedia, tais como imagens, referências, tabelas e listas.

# Passo a Passo da extração das árvores de análise sintática

## Passo 2:

Criação de um script em Java para automatizar a separação de sentenças pertencentes em um texto, onde o retorno da execução é um arquivo texto composto pelas sentenças extraídas do texto selecionado pelo usuário, separadas por linhas.

# Passo a Passo da extração das árvores de análise sintática

## Passo 3:

Passagem do arquivo retornado pelo algoritmo de separação de sentenças no script parser (algoritmo de análise sintática baseado no Graphia), que usa os pos-taggers mWANN-Tagger e o LX Tagger para “tagueamento” dos termos de cada sentença que serão formatados em forma de árvore e impressos ou no terminal ou em um arquivo texto.

# Algoritmo usado na análise das sentenças

Foi criado um script em java para facilitar a visualização das sentenças.

Pelo qual, é possível visualizar os textos originais e os em forma de árvore resultado dos algoritmo de extrações, consultar uma legenda que indica o significado das tags, e identificar se a sentença foi corretamente parseada ou não e corrigir as sentenças incorretas.

Após analisar todas as sentenças, é gerado um arquivo com o relatório da análise do texto. Além disso, também é possível gerar um gráfico comparativo entre as sentenças corretas e incorretas.

# Sistema de análise (Tela de análise)

Análise das árvores

Em 2009, a Rede Record comprou o documentário e passou a divulgar trechos do mesmo na emissora.

(CONJ de\_)

(N o)

)

(AP

(ADV poder)

(A de)

)

)

)

)

(PP

(P Marinho)

(NP

(N .)

)

)

)

)

Legenda:

"V": verbo

"N": substantivo

"ADP": adposição

"PRON": pronome

"ADJ"/"A": adjetivo

"DET": determinante

"PUNC"/"PNT": pontuação

"PRS": pronome (não possessivo)

"INTJ": interjeição

"NUM": número

"ADV": advérbio

"CJ"/"CONJ": conjunção

"CARD": cardinal

"ART": artigo

"P": preposição

0

Incorreto

Show Chart

Correto

62

Save State



# Sistema de análise (Tela de correção)

Correção de árvore

Seu empreendedorismo levou à constituição de um dos maiores império de comunicação do planeta e o fez figurar diversas vezes entre os homens mais ricos do mundo

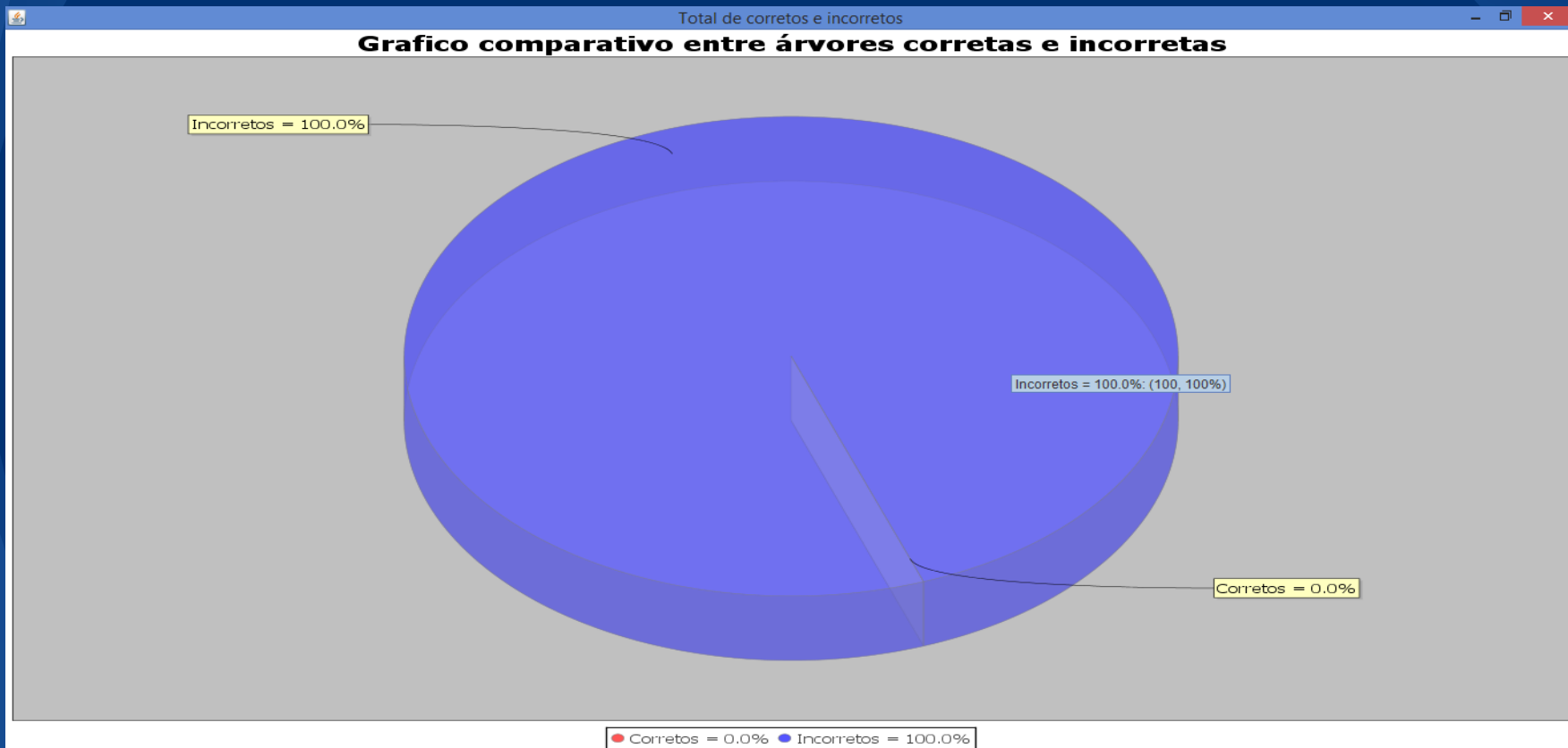
```
(ROOT
(S
(S
(S
(NP
(PRS Seu)
)
(VP
(V
(NP
(N empreendedorismo)
)
(V levou)
)
(ADV constituiu) //enconde //N
)
)
(NP
(CARD
(ADV
```

Legenda:

|                     |                                 |                        |
|---------------------|---------------------------------|------------------------|
| "V": verbo          | "DET": determinante             | "ADV": advérbio        |
| "N": substantivo    | "PUNC"/"PNT": pontuação         | "CJ"/"CONJ": conjunção |
| "ADP": adposição    | "PRS": pronome (não possessivo) | "CARD": cardinal       |
| "PRON": pronome     | "INTJ": interjeição             | "ART": artigo          |
| "ADJ"/"A": adjetivo | "NUM": número                   | "P": preposição        |

Salvar

# Resultados



# Resultados

Exemplo de saída do sistema de análise:

roberto marinho

Número de corretos:

0

0.0%

Número de incorretos:

69

100.0%

Sentenças corretas:

Sentenças incorretas:

Roberto Marinho

Roberto Pisani Marinho (Rio de Janeiro, — Rio de Janeiro, ) foi um jornalista e empresário brasileiro

Proprietário das Organizações Globo de 1925 a 2003, foi um dos homens mais poderosos e influentes do país no século XX

# Resultados

A análise dos textos trouxeram um resultado um pouco inesperado.

100% das sentenças estavam com algum tipo de incorreção. Salvo apenas por algumas sentenças formadas por títulos que possuem apenas 1 ou 2 palavras, entretanto, essas sentenças são desconsideradas para a análise.

# Problemas encontrados nos resultados

Algumas palavras estavam incompletas devido a presença de algum tipo de acentuação (provavelmente não houve um tratamento de codificação de palavras no algoritmo), e isso fazia com que houvesse uma classificação incorreta.

Além desse problema de encode, a palavra Marinho na entidade nomeada Roberto Marinho foi identificado como adjetivo que é uma das classificações possíveis para a palavra, apesar de estar errado no contexto.

Houveram também, algumas palavras que foram caracterizadas totalmente erradas.

# Dificuldades encontradas

- A separação de sentenças foi errada em casos que pontuação não indicassem final de frase, porém foi indiferente nessa fase da análise.
- Alguns bugs foram encontrados no parser, pois ele retornava um erro quando executava com o texto do Rio de Janeiro.
- Dificuldade para lembrar todas as regras gramaticais para realizar a análise.

# Trabalhos Futuros

Como trabalho futuro, possuímos a intenção de remodelar o algoritmo do Graphia para que funcione corretamente para textos em português, pois os resultados dessa análise foram negativos e a proposta desse trabalho era avaliar a necessidade de fazer essa remodelagem.

