

Processamento lexical, morfologia e morfossintaxe

SCC5908 Introdução ao
Processamento de Língua Natural

Thiago A. S. Pardo

Preâmbulo

- ▶ Em **processamento de texto**, é comum
 - Substituir uma palavra por outra
 - Procurar por uma informação, como data, nome, etc.
 - Analisar determinadas palavras
 - Mais genericamente, **procurar por padrões** no texto
 - Padrões simples: palavras
 - Padrões mais complexos: expressões, segmentos maiores

Exemplo

- ▶ Busca por todos os **valores monetários** em um texto

Levantamento da consultoria Economática aponta que empresas brasileiras de capital aberto tiveram os maiores lucros em 2010, considerando somente o setor de construção civil das Américas.

*No topo da lista está a PDG Realty, com um lucro líquido de **US\$ 473,9 milhões** apurados em 2010, bem acima da segunda colocada, a americana Walter Industries, com ganhos de **US\$ 385,8 milhões** para o mesmo período.*

As companhias brasileiras ocupam as próximas cinco posições (da 3ª à 7ª) no ranking preparado pela Economática, a saber: MRV, Cyrela, Gafisa, Brookfield e Rossi Residencial. Numa lista de dez posições, sete são ocupados por companhias nacionais.

*A consultoria também preparou uma lista baseada em volume de vendas, desta vez com presença mais expressiva das construtoras americanas, a exemplo do primeiro lugar -- a Pulte Homes-- com um volume de **US\$ 4,44 bilhões** em imóveis comercializados, e do segundo lugar --a Horton-- com uma cifra de **US\$ 3,96 bilhões**.*

...

3

Processamento textual

- ▶ **Útil** para
 - Tarefas particulares: buscar algo que leu
 - Tarefas científicas: sintomas e tratamentos de uma doença
 - Tarefas comerciais: sistemas on-line

4

Apple iPad 3G 64 GB - Ta... x




www.bondfaro.com.br/preco--tablet--apple-ipad-3g-64-gb.html


BONDfaro Faça parte da Comunidade Login | Cadastre-se

ipad **BUSCAR** Todas as categorias

Início > Informática > Tablet > Apple iPad 3G 64 GB

Apple iPad 3G 64 GB Categoria: Tablet

Geral Preços Fotos e vídeos Avaliação de quem comprou Compartilhar:   



Se você já usou este produto, que tal enviar sua avaliação e ajudar outras pessoas a decidirem sua compra?

ENVIE SUA AVALIAÇÃO

Nota dos Usuários

★★★★★
Baseado em 1 opinião

[Leia todas as avaliações >>](#)

Resumo das especificações

R\$ 2.031,42 - R\$ 2.399,00 (em 10 lojas abaixo)

Marca	Apple
Modelo	iPad 3G
Conexões	3G, Bluetooth, Fone de Ouvido, USB, Wi-Fi
Formatos Aceitos	AAC, DOC, GIF, JPEG, MP3, MP4, PDF, PPT, WAV, WMA

[Veja a Ficha Técnica completa >>](#)







Comunidade

2 fãs 0 manual 21 fotos 7 Vídeos

Apple iPad 3G 64 GB - Ta... x

www.bondfaro.com.br/preco--tablet--apple-ipad-3g-64-gb.html

Onde comprar (10 lojas) Popularidade **Menor Preço** Loja

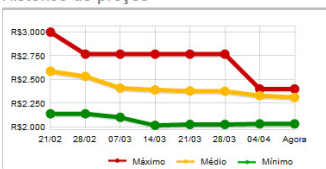
	R\$ 2.399,00	12x de R\$ 199,92	IR À LOJA
	R\$ 2.159,10		IR À LOJA
	R\$ 2.031,42	12x de R\$ 199,16	IR À LOJA
	R\$ 2.399,00	12x de R\$ 199,92	IR À LOJA
	R\$ 2.399,00	12x de R\$ 199,92	IR À LOJA
	R\$ 2.159,10	12x de R\$ 199,92	IR À LOJA

BuscaPé

No BuscaPé você economiza e concorre a prêmios

Participe

Histórico de preços



Alerta de preços

Me avisar por e-mail:

Toda vez que o menor preço do site mudar

Quando o produto atingir um preço abaixo de R\$

The screenshot shows a search for 'ipad' on the BuscaPé website. The search results are categorized under 'Tablet'. Two main results are visible:

Product	Lojas	Ofertas	Price
Apple iPad 3G 64 GB	38 Lojas	41 Ofertas	R\$ 1.899,90 até R\$ 3.229,00
Apple iPad 3G 32 GB	33 Lojas	36 Ofertas	R\$ 1.673,07 até R\$ 2.749,99

Each product listing includes a star rating, a 'Ver todas' link, and a 'Compare Preços' button. The website header includes navigation links like 'Registre-se', 'Guia para Comprar pela Internet', and 'Minha conta'. A sidebar on the left lists various categories like 'Tablet', 'Acessórios para Tablet', 'Livros', etc.

The screenshot shows the WolframAlpha website with the search input 'dollar'. The results display the current exchange rate and a historical chart.

Assuming "dollar" is a unit | Use as a character or a word instead

Input interpretation:
\$1 (US dollar)

Local currency conversion for 1 CalculateScan MoneyScanner Private unit\$606008:
R\$1.60 (Brazilian reais) (at current quoted rate)

Exchange history for \$1 (US dollar):

Period	Value	Date
1-year minimum	R\$1.61	(05.04.2011 1 day ago)
1-year maximum	R\$1.89	(20.05.2010 11 months ago)

The chart shows the exchange rate fluctuating between approximately R\$1.60 and R\$1.89 over the period shown. A 'Units' link is visible at the bottom right.

Exemplo

- ▶ Busca por todos os **valores monetários** em um texto

Levantamento da consultoria Economática aponta que empresas brasileiras de capital aberto tiveram os maiores lucros em 2010, considerando somente o setor de construção civil das Américas.

*No topo da lista está a PDG Realty, com um lucro líquido de **US\$ 473,9 milhões** apurados em 2010, bem acima da segunda colocada, a americana Walter Industries, com ganhos de **US\$ 385,8 milhões** para o mesmo período.*

As companhias brasileiras ocupam as próximas cinco posições (da 3ª à 7ª) no ranking preparado pela Economática, a saber: MRV, Cyrela, Gafisa, Brookfield e Rossi Residencial. Numa lista de dez posições, sete são ocupados por companhias nacionais.

*A consultoria também preparou uma lista baseada em volume de vendas, desta vez com presença mais expressiva das construtoras americanas, a exemplo do primeiro lugar -- a Pulte Homes-- com um volume de **US\$ 4,44 bilhões** em imóveis comercializados, e do segundo lugar --a Horton-- com uma cifra de **US\$ 3,96 bilhões**.*

...

Como
fariam?

9

Expressões Regulares (ER)

- ▶ Notação tradicional para caracterizar segmentos textuais de todo tipo
 - Especificam **seqüências de símbolos** a serem buscados/caracterizados
 - Vários sistemas de busca de expressões regulares
 - grep, no Linux/UNIX
 - Lex/flex
 - Há variações de sistema para sistema, mas são muito parecidas

10

ER: notação

▶ Exemplos

- Casamento direto: **preço**
- Letra maiúscula ou minúscula no início: **[Pp]reço**
 - [] indicam disjunção, ou seja, um único elemento do conjunto
- Identificação de um único dígito do texto: **[0123456789]**
- Identificação de uma letra em um intervalo de letras: **[a-z]**
- Qualquer caractere diferente de a: **^a**

11

ER: notação

▶ Exemplos

- Singular ou plural: **preços?**
- 1 ou mais ocorrências (+) de algum elemento: **Aa+i!**
 - Aai!, Aaaaaaiiii!
- 0 ou mais ocorrências (*) de algum elemento: **Aa*i*!**
 - Aaaaiii!, Aaiiii!, Ai!, Aaaa!
- Caractere curinga (.): **beg.n**
 - begin, began, begun
- Alternativa (|): **preço|os** ou **(gato)|(cão)**
 - O que acontece se tivermos **gato|cão** sem parênteses?

12

Exercícios

- ▶ Como identificar nomes próprios?
- ▶ E e-mails?

13

Exercícios

- ▶ Como identificar nomes próprios?
 - `[A-Z][a-z]+`
 - ▶ E e-mails?
 - `[a-z0-9_]+@[a-z\.]+`
- **Cuidado:** alguns caracteres são especiais e, para serem usados em seu sentido original, precisam de `\` ou `""`
 - Exemplos: `.` `$` `-`

14

Exercício

- ▶ Expressão regular para reconhecer os valores monetários?

Levantamento da consultoria Economática aponta que empresas brasileiras de capital aberto tiveram os maiores lucros em 2010, considerando somente o setor de construção civil das Américas.

*No topo da lista está a PDG Realty, com um lucro líquido de **US\$ 473,9 milhões** apurados em 2010, bem acima da segunda colocada, a americana Walter Industries, com ganhos de **US\$ 385,8 milhões** para o mesmo período.*

As companhias brasileiras ocupam as próximas cinco posições (da 3ª à 7ª) no ranking preparado pela Economática, a saber: MRV, Cyrela, Gafisa, Brookfield e Rossi Residencial. Numa lista de dez posições, sete são ocupados por companhias nacionais.

*A consultoria também preparou uma lista baseada em volume de vendas, desta vez com presença mais expressiva das construtoras americanas, a exemplo do primeiro lugar -- a Pulte Homes-- com um volume de **US\$ 4,44 bilhões** em imóveis comercializados, e do segundo lugar --a Horton-- com uma cifra de **US\$ 3,96 bilhões**.*

...

15

Exercício

- ▶ Expressão regular para reconhecer os valores monetários?

Levantamento da consultoria Economática aponta que empresas brasileiras de capital aberto tiveram os maiores lucros em 2010, considerando somente o setor de construção civil das Américas.

*No topo da lista está a PDG Realty, com um lucro líquido de **US\$ 473,9 milhões** apurados em 2010, bem acima da segunda colocada, a americana Walter Industries, com ganhos de **US\$ 385,8 milhões** para o mesmo período.*

As companhias brasileiras ocupam as próximas cinco posições (da 3ª à 7ª) no ranking preparado pela Economática, a saber: MRV, Cyrela, Gafisa, Brookfield e Rossi Residencial. Numa lista de dez posições, sete são ocupados por companhias nacionais.

*A consultoria também preparou uma lista baseada em volume de vendas, desta vez com presença mais expressiva das construtoras americanas, a exemplo do primeiro lugar -- a Pulte Homes-- com um volume de **US\$ 4,44 bilhões** em imóveis comercializados, e do segundo lugar --a Horton-- com uma cifra de **US\$ 3,96 bilhões**.*

...

US\\$ [0-9]+,[0-9]+ [mb]ilhões

16

Exercício

Sentença:

O *homem* viu a mulher de binóculos na montanha.

Análise automática:

*O*_DET *homem*_N viu_V a_DET *mulher*_N de_PRP *binóculos*_N
em_PRP a_DET *montanha*_N.

Expressão regular para os **substantivos** e os **verbos**?

17

Exercício

Sentença:

O *homem* viu a mulher de binóculos na montanha.

Análise automática:

*O*_DET *homem*_N viu_V a_DET *mulher*_N de_PRP *binóculos*_N
em_PRP a_DET *montanha*_N.

Expressão regular para os **substantivos** e os **verbos**?

[A-Za-z][a-z]*_N|V

18

Exercício

Sentença:

O *homem* viu a mulher de binóculos na montanha.

Análise automática:

*O_DET homem_N viu_V a_DET mulher_N de_PRP binóculos_N
em_PRP a_DET montanha_N.*

Expressão para **substantivos seguidos de verbos?**

19

Exercício

Sentença:

O *homem* viu a mulher de binóculos na montanha.

Análise automática:

*O_DET homem_N viu_V a_DET mulher_N de_PRP binóculos_N
em_PRP a_DET montanha_N.*

Expressão para **substantivos seguidos de verbos?**

[A-Za-z][a-z]*_N [a-z]+_V

20

Autômatos

- ▶ **Expressões regulares** implementadas como **autômatos de estados finitos**
 - Autômato: modelo matemático eficaz e elegante para lidar com expressões regulares
- ▶ Autômatos utilizados para revisão ortográfica, síntese e reconhecimento de fala, extração de informação, tradução automática, **análise morfológica**, análise morfossintática, etc.

21

Autômatos

- ▶ Poder representacional equivalente



22

Autômatos

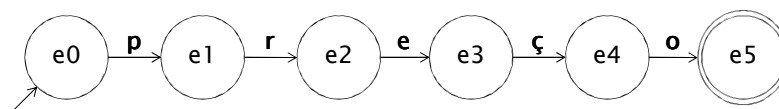
▶ Componentes

- **Estados** que modelam o “sistema”
 - Pontos da análise sendo realizada, por exemplo
- **Símbolos de entrada**
 - Letras das palavras, números, símbolos, etc.
- **Estados inicial e final**
 - Início e fim do processo
- **Transições** entre estados

23

Exemplo

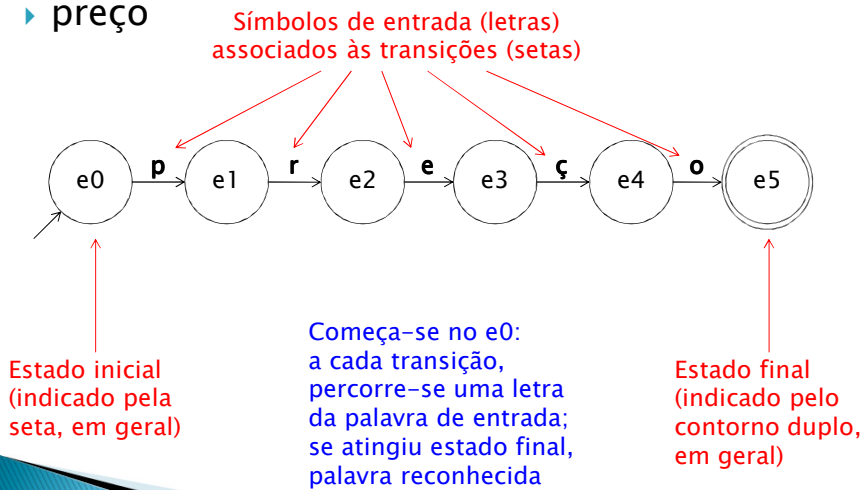
▶ preço



24

Exemplo

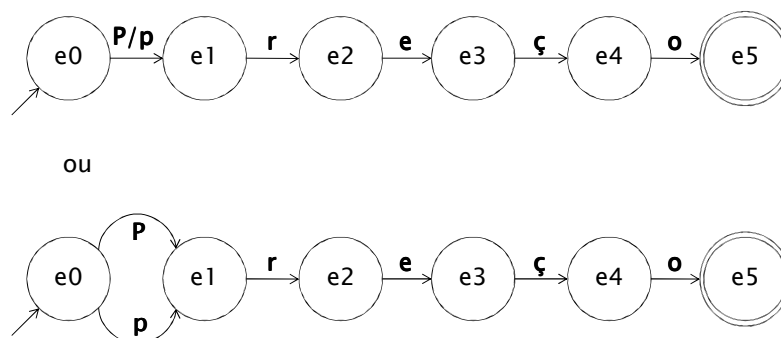
► preço



25

Exemplo

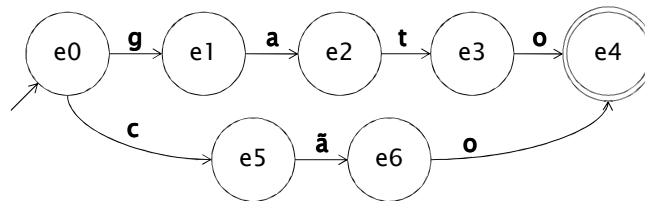
► [Pp]reço



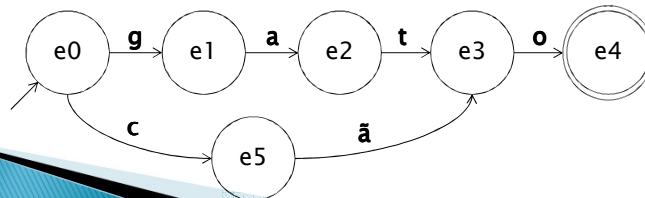
26

Exemplo

► (gato)|(cão)



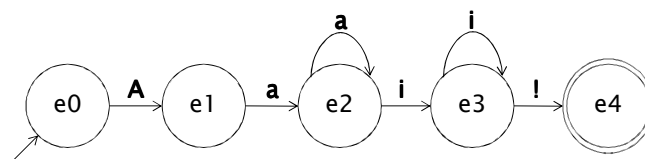
ou



27

Exemplo

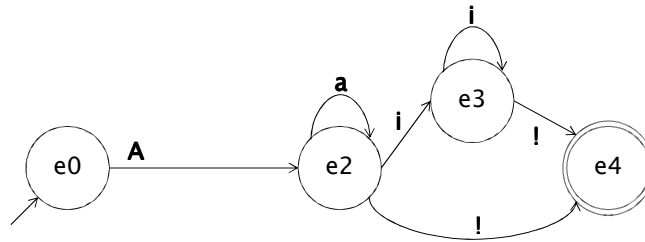
► Aa+i+!



28

Exemplo

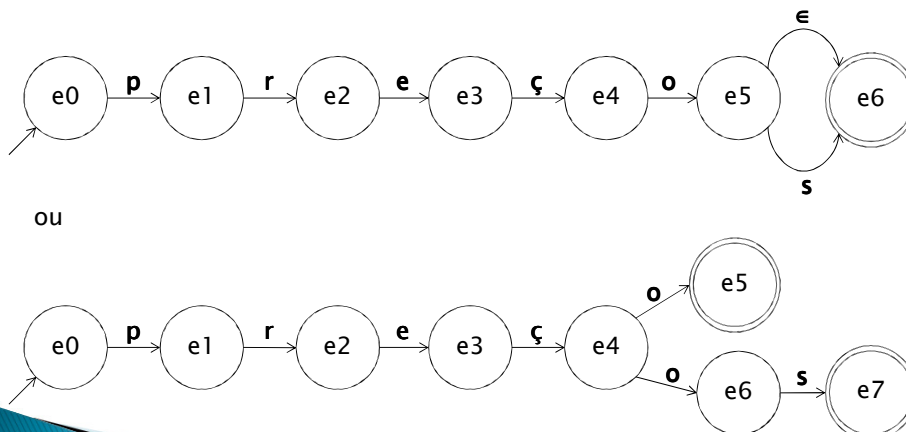
▶ Aa^*i^*



29

Exemplo

▶ preços?



30

Exercício

- ▶ Criar autômato para reconhecer **valores monetários**
 - US\\$ [0-9]+,[0-9]+ [mb]ilhões

31

Autômatos

- ▶ **Variações**
 - Transdutores
 - Além de reconhecerem a entrada, geram saída
 - Usados em análise morfológica
 - Modelos de Markov
 - Redes de transição

32

Análise morfológica

▶ *Parsing* morfológico

- Analisar uma palavra e identificar seus componentes
 - Morfemas
 - Possibilidades
 - meninos → lema (menino), masculino (o), plural (+s), subst
 - meninos → radical (menin), masculino (+o), plural (+s), subst
 - meninas → lema (menino), feminino (-o +a), plural (s), subst

33

Análise morfológica

▶ Relevância da tarefa

- ???

34

Análise morfológica

▶ Relevância da tarefa

- Reconhecer palavras e suas variações
 - Revisão ortográfica, busca na web, sumarização, extração de informação
 - Stemming, lematização
- ... e também produzir a forma adequada das palavras, derivar palavras novas, lidar com neologismos
 - Geração textual, tradução automática
 - “Máquina morfológica”
- Caracterização léxica da língua, no geral

35

Terminologia básica

▶ **Morfemas**: unidade mínima de significado

- Raiz/radical
 - Alguns diferenciam esses termos, outros não
- Afixos

▶ **Afixos**

- Prefixo: desamor, infeliz
- Sufixo: lealdade, facilmente, quebrado, comia
- Infixo: rabiscar
 - Raro, alguns dizem que não existe para o português
- Circunfixo: anoitecer, descampado

36

Terminologia básica

▶ Morfe

- Realização de um morfema
 - Morfema é abstrato, enquanto morfe é concreto
 - Exemplo: morfema de negação pode ser expresso pelos morfes in (de infeliz) ou i (de imutável)

▶ Alomorfes

- Morfes que expressam um mesmo morfema
 - In e i para negação
 - Ante, pré e pró para anterioridade

37

Terminologia básica

▶ Processos principais de formação de palavras

- **Flexional**: variações de uma mesma palavra
 - Flexão nominal: número, gênero
 - Flexão verbal: modo-tempo, número-pessoa
 - Adição de morfemas gramaticais
- **Derivacional**: palavras novas
 - Podem mudar classe e sentido
 - “modelo” → “modelagem”
 - Adição de morfemas lexicais

38

Análise morfológica

- ▶ Para **construir um parser morfológico**, são necessários
 - **Léxico**
 - Radicais e afixos e suas possíveis classificações (substantivos, verbos, etc.)
 - Conhecimento de **morfotática**
 - Como os morfemas se ordenam para que as palavras se formem
 - Exemplo: em português, o morfema de plural aparece após e substantivo, e não antes
 - “Sintaxe da morfologia”
 - **Regras ortográficas**
 - Modelam mudanças que ocorrem nas palavras quando morfemas se combinam
 - Exemplo: casa+PL=casaS, mas flor+PL=florES

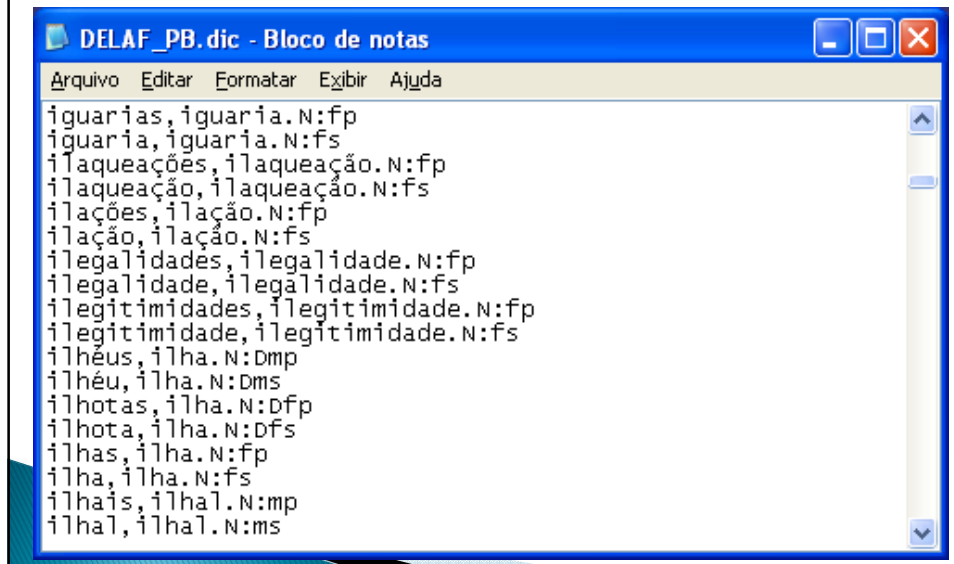
39

Análise morfológica

- ▶ **Alternativa 1**
 - Listagem de palavras
 - Exaustiva: léxico de formas analisadas (também chamadas flexionadas ou plenas)
 - Palavras com todas as suas variações
 - Pouca economia, redundância, compactação de arquivos

40

Exemplo do UNITEX-PB



Análise morfológica

► Alternativa 1

- Listagem de palavras
 - Econômica: léxico de raízes (ou de morfemas)
 - Listagem de raízes + regras de formação das palavras (morfológica e regras ortográficas)
 - Mais economia, processo mais caro

Análise morfológica

▶ Alternativa 1

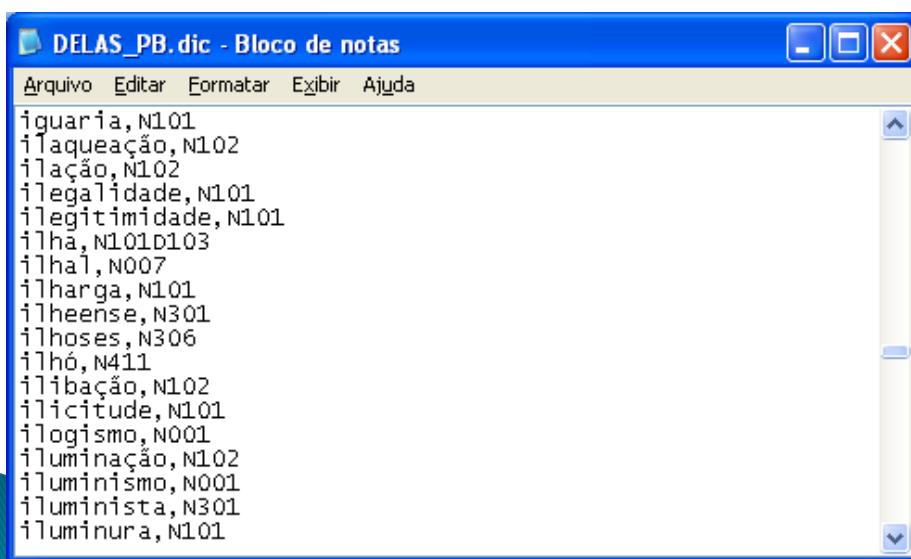
◦ Listagem de palavras

• Meio termo

- Léxico de lemas (ou formas canônicas) associados as suas variações
- Palavras irregulares em formas plenas no léxico + léxico de raízes para palavras regulares
- Etc.

43

Exemplo do UNITEX-PB



```
DELAS_PB.dic - Bloco de notas
Arquivo  Editar  Formatar  Exibir  Ajuda
íguaríá, N101
ílaqueação, N102
ílação, N102
ílegalidade, N101
ílegitimidade, N101
ílhá, N101D103
ílhá1, N007
ílharga, N101
ílheense, N301
ílhoses, N306
ílhó, N411
ílibação, N102
ílicitude, N101
ílogismo, N001
íluminação, N102
íluminismo, N001
íluminista, N301
íluminura, N101
```

Análise morfológica

► Alternativa 1

◦ Listagem de palavras

- **Problemas** para lidar com
 - Novas palavras e variações: novos verbos (denominais, inclusive), nomes próprios, etc.
 - Línguas morfológicamente complexas
 - Turco, por exemplo

45

Turco – exemplo

(Jurafsky e Martin, 2008)

uygarlaştiramadıklarımızdanmışsınızcasına

uygar +laş +tır +ama +dık +lar +ımız +dan +mış +sınız +casına
civilized +BEC +CAUS +NABL +PART +PL +P1PL +ABL +PAST +2PL +AsIf

“(behaving) as if you are among those whom we could not civilize”

+BEC	“become”
+CAUS	the causative verb marker (‘cause to X’)
+NABL	“not able”
+PART	past participle form
+P1PL	1st person pl possessive agreement
+2PL	2nd person pl
+ABL	ablative (from/among) case marker
+AsIf	derivationally forms an adverb from a finite verb

46

Análise morfológica

▶ Alternativa 2

- Codificação em forma de autômatos: maior eficiência computacional
 - De forma **complementar com o léxico**
 - Formas básicas/raízes no léxico e regras de formação de palavras (morfotática e regras ortográficas) mapeadas em autômatos
 - De forma **isolada**
 - Todo o léxico da língua mapeado em autômatos

47

Análise morfológica

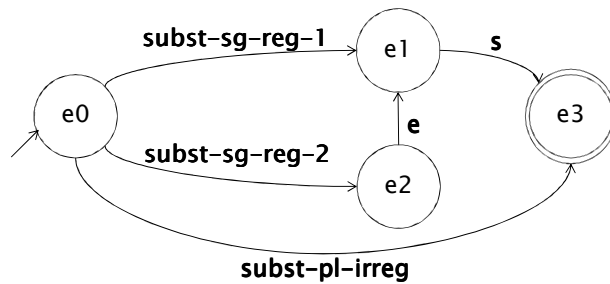
▶ Mapear palavras em seus componentes

- gatos → gato + SUBST + MASC + PL
- canto → canto + SUBST + MASC + SG
- canto → cantar + V + 1P + SG + Pind
 - A rigor, a tarefa de desambiguar “canto” (SUBST ou V) está além da análise morfológica
 - Morfossintaxe
 - No momento, listam-se todas as possibilidades

48

Exemplo simples

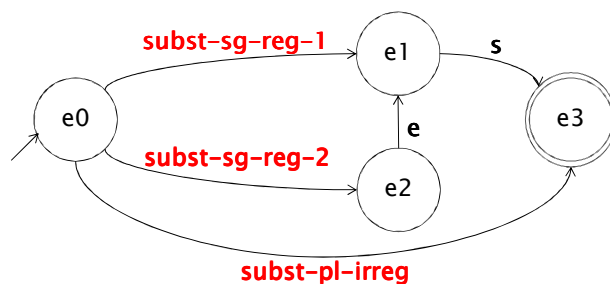
- ▶ Reconhecimento/geração de alguns substantivos no plural
 - Léxico de lemas + autômato



subst-sg-reg-1	subst-sg-reg-2	subst-pl-irreg
casa	flor	lápiz
porta	lar	córpuz
...

Exemplo simples

- ▶ Reconhecimento/geração de alguns substantivos no plural
 - Léxico de lemas + autômato



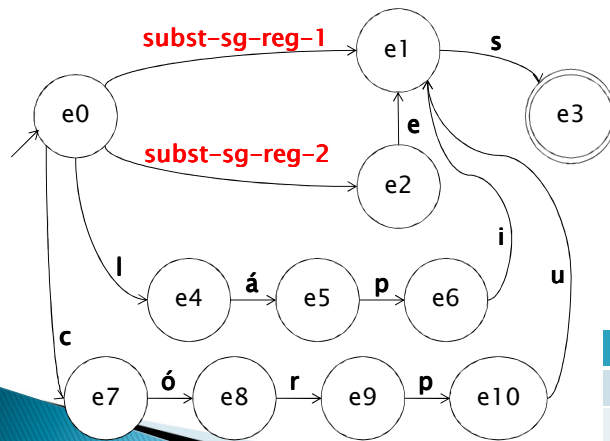
Podem ser substituídos pelos autômatos correspondentes!

Como?

subst-sg-reg-1	subst-sg-reg-2	subst-pl-irreg
casa	flor	lápiz
porta	lar	córpuz
...

Exemplo simples

- ▶ Reconhecimento/geração de alguns substantivos no plural
 - Léxico de lemas + autômato



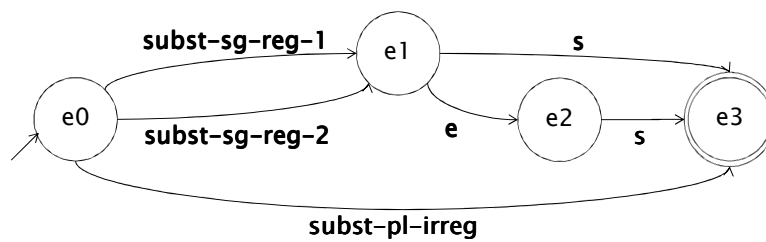
Podem ser substituídos pelos autômatos correspondentes!

Continuem!

subst-sg-reg-1	subst-sg-reg-2
casa	flor
porta	lar
...	...

Exemplo simples

- ▶ Cuidado com **overgeneration!**
 - O que acontece no caso abaixo?

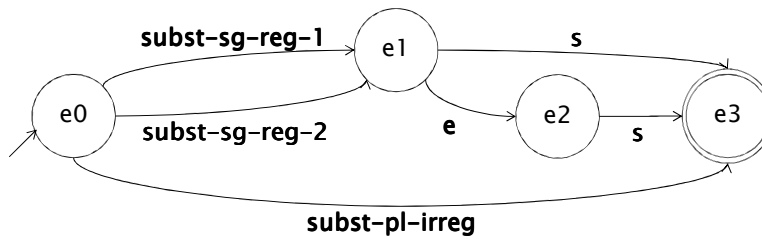


subst-sg-reg-1	subst-sg-reg-2	subst-pl-irreg
casa	flor	lápiz
porta	lar	cópus
...

Exemplo simples

- ▶ Cuidado com *overgeneration*!
- O que acontece no caso abaixo?

casas
*casaes
*flors
flores
...



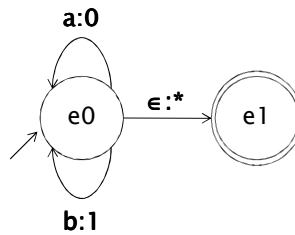
subst-sg-reg-1	subst-sg-reg-2	subst-pl-irreg
casa	flor	lápiz
porta	lar	cópus
...

Análise morfológica

- ▶ Para nossa tarefa, precisamos de mais poder
 - Além de se reconhecer/gerar as palavras, é necessário identificar os componentes
 - gatos → gato + SUBST + MASC + PL
 - canto → canto + SUBST + MASC + SG
 - canto → cantar + V + 1P + SG + Pind
 - Transdutores
 - Reconhecem a entrada e, em paralelo, geram saída

Transdutores

- ▶ Lendo a_s e b_s e gerando 0_s e 1_s , respectivamente, terminando com *

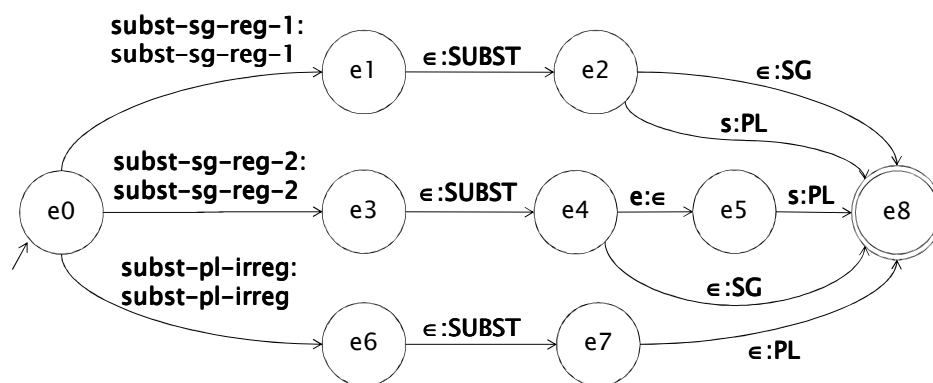


Análise de abba

55

Transdutores: exemplo

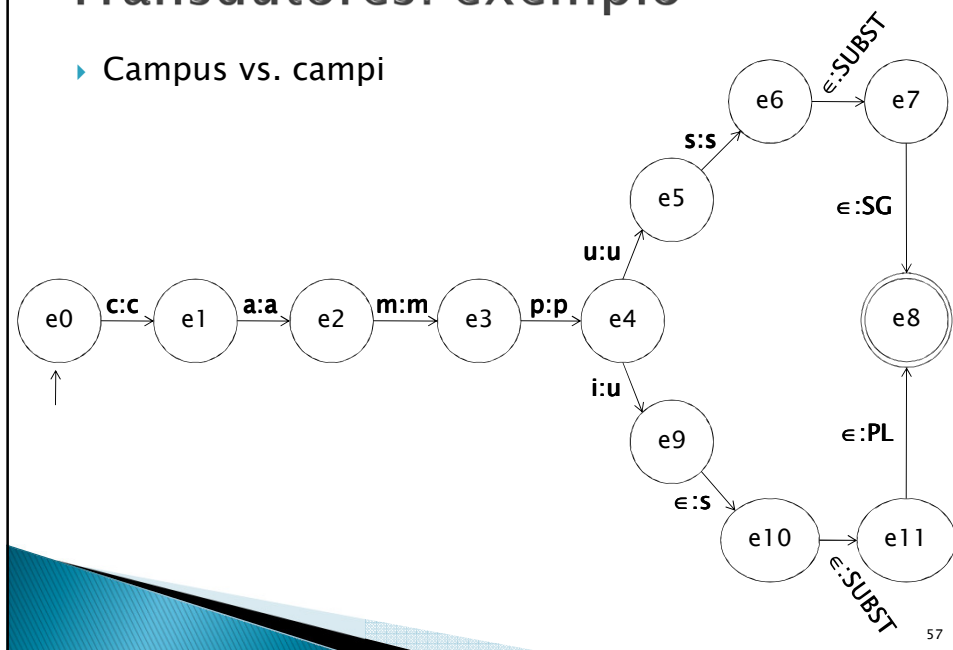
- ▶ Releitura do autômato de substantivos



subst-sg-reg-1	subst-sg-reg-2	subst-pl-irreg
casa	flor	lápiz
...

Transdutores: exemplo

- ▶ Campus vs. campi



Transdutores: exemplo

- ▶ **Menino, menina, meninos, meninas: exercício**
 - Reconhecer número, gênero, raiz e etiqueta morfossintática

Transdutores

- ▶ E casos como o de “canto”?
- Como identificar que “canto” pode ser um **verbo** ou um **substantivo**, gerando-se os atributos correspondentes para cada caso?
 - canto → canto + SUBST + MASC + SG
 - canto → cantar + V + 1P + SG + Pind

59

Transdutores

- ▶ E casos como o de “canto”?
- Como identificar que “canto” pode ser um **verbo** ou um **substantivo**, gerando-se os atributos correspondentes para cada caso?
 - canto → canto + SUBST + MASC + SG
 - canto → cantar + V + 1P + SG + Pind
- **A palavra seria reconhecida por mais de um transdutor!**
 - **Análise morfossintática** para desambiguar

60

Origens da Morfossintaxe

- ▶ **Dionísio Trácio, 100 AC**
 - Esboço da gramática do grego
 - Cunhou o vocabulário atual
 - Sintaxe, ditongo, clítico, etc.
 - 8 etiquetas morfossintáticas: substantivo, verbo, pronome, preposição, advérbio, conjunção, particípio, artigo
 - Vocabulário usado até hoje!
- ▶ **Morfossintaxe**
 - Morfologia: tipos de afixos possíveis variam com a classe
 - Sintaxe: palavras com comportamentos/funções similares em seus contextos são de uma mesma classe
 - Algo mais?

61

Origens da Morfossintaxe

- ▶ **Dionísio Trácio, 100 AC**
 - Esboço da gramática do grego
 - Cunhou o vocabulário atual
 - Sintaxe, ditongo, clítico, etc.
 - 8 etiquetas morfossintáticas: substantivo, verbo, pronome, preposição, advérbio, conjunção, particípio, artigo
 - Vocabulário usado até hoje!
- ▶ **Morfossintaxe**
 - Morfologia: tipos de afixos possíveis variam com a classe
 - Sintaxe: palavras com comportamentos/funções similares em seus contextos são de uma mesma classe
 - **Semântica: substantivos têm uma preferência por objetos, lugares e coisas, adjetivos por propriedades, etc.**
 - **Pragmática**

62

Conjuntos de etiquetas

► Variam muito

- Penn Treebank (Marcus et al., 1993): 45
- Brown Corpus (Francis, 1979): 87
- CLAWS 7 (Garside et al. 1997): 146
- Palavras (Bick, 2000): 14
- Mac-Morpho/Lácio-Web (Aluísio et al., 2003): 31

63

Exemplo: Penn Treebank

CC	Coordinating conjunction		
CD	Cardinal number		
DT	Determiner		
EX	Existential <i>there</i>		
FW	Foreign word		
IN	Preposition or subordinating conjunction		
JJ	Adjective		
JJR	Adjective, comparative		
JJS	Adjective, superlative		
LS	List item marker		
MD	Modal		
NN	Noun, singular or mass		
NNS	Noun, plural		
NP	Proper noun, singular	RB	Adverb
NPS	Proper noun, plural	RBR	Adverb, comparative
PDT	Predeterminer	RBS	Adverb, superlative
POS	Possessive ending	RP	Particle
PP	Personal pronoun	SYM	Symbol
PP\$	Possessive pronoun	TO	to
		UH	Interjection
		VB	Verb, base form
		VBD	Verb, past tense
		VEG	Verb, gerund or present participle
		VBN	Verb, past participle
		VBP	Verb, non-3rd person singular present
		VBZ	Verb, 3rd person singular present
		WDT	Wh-determiner
		WP	Wh-pronoun
		WP\$	Possessive wh-pronoun
		WRB	Wh-adverb

Exemplo: Mac-Morpho

Tag	Definition	Compl. Tag	Definition
ADJ	open-class noun modifier	EST	foreign
ADV-KS-REL	relative subordinating Adverb	AP	apposition
ADV-KS	Non-relative subordinating Adverb	+	contraction/ enclitic
ADV	Non-subordinating adverb	!	mesoclitic
ART	Article	[beginning,
KC	coordinating conjunction	...	middle part,
KS	coordinating conjunction]	and end of discontinuous compound (further discussed in Section 3)
IN	interjection	TEL	phone number
N	open-class noun phrase nucleus	DAT	date
NPROP	proper noun	HOR	time
NUM	numeral as a noun modifier	DAD	formatted data not falling into above categories
PCP	past participle or adjective		
PDEN	emphasis/focus		
PREP	preposition		
PROPESS	personal pronoun		
PRO-KS-REL	relative subordinating pronoun		
PRO-KS	Non-relative subordinating pronoun		
PROSUB	non-subordinating pronoun as a noun phrase nucleus		
PROADJ	Non-subordinating pronoun as a modifier		
VAUX	Auxiliary verb		
V	Non-auxiliary verb		
CUR	Currency symbol		

65

Terminologia

- ▶ **Palavras de classes fechadas, palavras funcionais**
 - Conjunto de palavras varia pouco
 - Preposições, conjunções, artigos
 - ▶ **Palavras de classes abertas**
 - O conjunto varia bastante, surgindo novas palavras
 - Substantivos, verbos
- Conjuntos de palavras de classes abertas e fechadas não são iguais para todas as línguas
- Nem todas as classes existem para todas as línguas ou são distinguidas das mesmas formas

66

Terminologia

- ▶ Substantivos/nomes
 - Comuns, próprios
 - Contáveis (abelha, casa), incontáveis (ar, açúcar)
- ▶ Verbos
 - Principais, auxiliares
- ▶ Advérbios
 - Tempo, local, modo, direção, etc.
- ▶ Conjunções
 - Coordenativas e subordinativas
- ▶ Pronomes
 - Pessoais, possessivos, interrogativos, relativos, etc.

67

Etiquetas morfossintáticas

- ▶ **Nem sempre a distinção é simples**
 - Advérbios vs. preposições
 - *Near, around*
 - Adjetivos vs. participípios
 - Eles estão casados.
 - Advérbios: tudo que não cabe nas outras classes

68

Etiquetação morfossintática

- ▶ **Tagging**, ou **parsing morfossintático**
 - Associação de etiquetas às palavras de uma sentença
 - Faz-se necessário, portanto, tokenização e segmentação sentencial
 - Tarefa de desambiguação: dentre as etiquetas (tags) possíveis previstas (pelo léxico, por exemplo), determinar a mais apropriada
 - **Contexto desambigua!**

69

Tagging

- ▶ **Útil** para um infinidade de tarefas de PLN
 - *Stemming*, lematização
 - Tradução, sumarização, auxílio à escrita
 - Identificação de autoria, extração de informação
 - Pesquisas lingüísticas variadas: neologismos, comportamento de palavras, etc.
 - Etc.

70

Tagging

▶ 2 principais abordagens

◦ Regras

- Por exemplo, uma palavra antecida por um artigo é um substantivo

◦ Probabilidades

- Classe mais provável de uma palavra em função das palavras vizinhas, com aprendizado a partir de cópuz

▶ Hibridismo também é possível

- Por exemplo, aprendizado de regras a partir de cópuz

71

Tagging: regras

▶ Primeiras abordagens (década de 60)

◦ 2 passos tradicionais

- Léxico fornece possíveis classes para cada palavra
- Regras criadas manualmente são utilizadas para desambiguar

▶ Mais recentemente

- Dicionários maiores e muito mais regras!

72

Exemplo: EngCG tagger (Voutilainen, 1999)

→ Análise morfológica da sentença (tag correta em negrito)

Pavlov had shown that salivation.

Pavlov:	PAVLOV N NOM SG PROPER
had:	HAVE V PAST VFIN SVO HAVE PCP2 SVO
shown:	SHOW PCP2 SVOO SVO SV
that:	ADV PRON DEM SG DET CENTRAL DEM SG CS
salivation:	N NOM SG
::	PUNC DOT



73

Exemplo: EngCG tagger (Voutilainen, 1999)

→ Aplicação de regras para determinar as melhores tags

Exemplo de regra

WORD: that

IF

next word is adj, adverb, or quantifier AND

after this word there is the sentence boundary AND

the previous word is not a verb that allows adjs as complements

THEN eliminate non-adv tags

ELSE eliminate adv tag

74

Exemplo: ReGra (Martins et al., 1998)

Exemplo de entrada (com erros) para o revisor gramatical

OS	Definite article (the): masculine, plural Personal pronoun (them): masculine, plural
MENINO	Noun (boy): masculine, singular
PREFERE	Verb (to prefer): 3rd person, singular, present tense, indicative, transitive
BRINCAR	Verb (to play): infinitive
DO	Contraction: preposition (of) + definite article (the): masculine, singular
QUE	Relative Pronoun (which) Adverb (what) Conjunction (than)
ESTUDAR	Verb (to study): infinitive

Regras de desambiguação utilizadas

OS	Definite article (the): the following word is a masculine noun
QUE	Conjunction (than): the previous word is a contraction (preposition + article)

75

Tagging: regras

- ▶ **Zellig Harris** (1962) e o primeiro tagger (provavelmente)
 - 14 regras de desambiguação
- ▶ **UNITEX-PB** (Muniz, 2004)
 - 80 regras de desambiguação no formalismo ELAG

ELAG
16 regras para Adjetivos
30 regras para Advérbios
22 regras para Artigos
12 regras para Substantivos

76

Tagging: regras

- ▶ ELAG (*Elimination of Lexical Ambiguities by Grammars*) (Laporte e Monceaux, 1998)

As seguintes premissas são seguidas pelo formalismo ELAG: análises corretas não devem ser removidas; os **resultados de análise sintática não podem ser explicitamente utilizados**, uma vez que eles não estão disponíveis quando a resolução de ambigüidade lexical é aplicada ao texto; a análise lingüística que desejamos aplicar à sentença deve ser levada em consideração, o que implica que **o criador das gramáticas de resolução de ambigüidade lexical tem visões particulares** sobre o resultado desejado da análise sintática.

Muniz (2004)