

# 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 Pad 3G 64 GB

Apple iPad 3G 64 GB Categoria: Tablet

Geral Preços Fotos e vídeos Avaliação de quem comprou Compartilhar:   



Passa o mouse sobre a imagem para ver mais detalhes



**Nota dos Usuários**

★★★★★  
Baseado em 1 opinião

[Leia todas as avaliações >>](#)

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

[ENVIE SUA AVALIAÇÃO](#)

**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

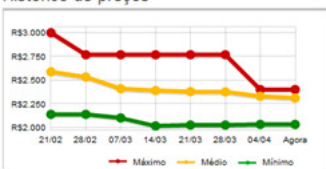
	R\$ 2.399,00	12x de R\$ 199,92	<a href="#">IR À LOJA</a>
	R\$ 2.159,10		<a href="#">IR À LOJA</a>
	R\$ 2.031,42	12x de R\$ 199,16	<a href="#">IR À LOJA</a>
	R\$ 2.399,00	12x de R\$ 199,92	<a href="#">IR À LOJA</a>
	R\$ 2.399,00	12x de R\$ 199,92	<a href="#">IR À LOJA</a>
	R\$ 2.159,10	12x de R\$ 199,92	<a href="#">IR À LOJA</a>

**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\$ 2.031,42

The screenshot shows a search engine interface for 'BuscaPé'. The search term 'ipad' is entered in the search bar. The results page features a sidebar with categories like 'Tablet', 'Acessórios para Tablet', 'Livros', 'Capa para Celular e Smartphone', and 'Carregador para Celular e Smartphone'. The main content area displays two product listings for the Apple iPad 3G. The first listing is for the 64 GB model, priced at R\$ 1.899,90, with 38 stores and 41 offers. The second listing is for the 32 GB model, priced at R\$ 1.673,07, with 33 stores and 36 offers. Both listings include star ratings and a 'Compare Preços' button.

The screenshot shows the WolframAlpha website with the search term 'dollar'. The results display the conversion of 1 US dollar to Brazilian Reals. The current rate is R\$1.60. Below this, there is an 'Exchange history for \$1 (US dollar)' section with a line graph showing the rate from July to April. The 1-year minimum is R\$1.61 (05.04.2011 | 1 day ago) and the 1-year maximum is R\$1.89 (20.05.2010 | 11 months ago). The page number '8' is visible in the bottom right corner.

## 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\*!**
  - Aaaaii!, 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

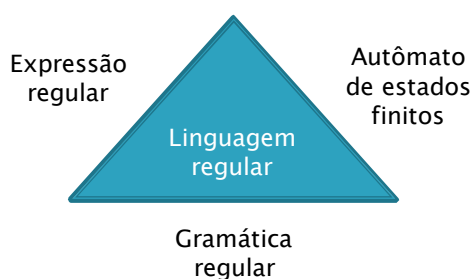
## Autômatos

- ▶ **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

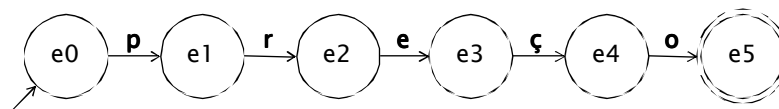
## ▶ Componentes

- **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

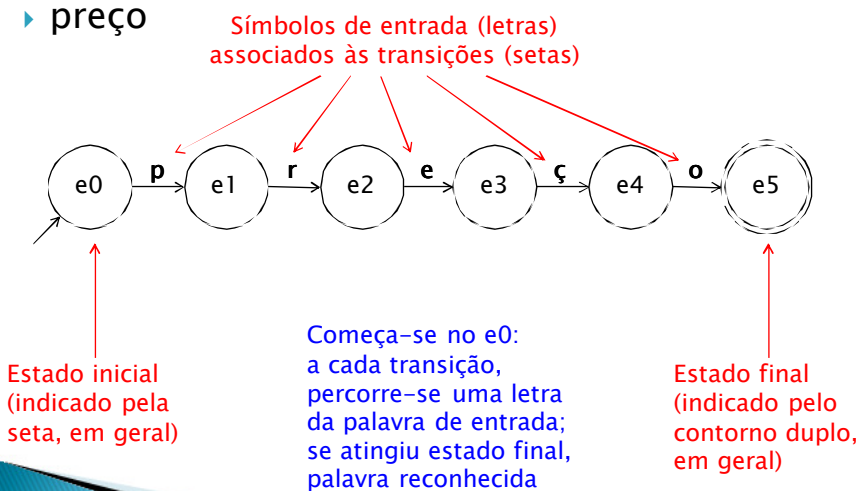
## ▶ preço



24

## Exemplo

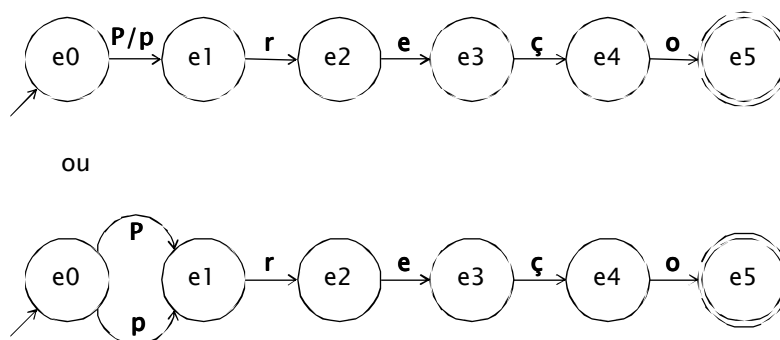
### ► preço



25

## Exemplo

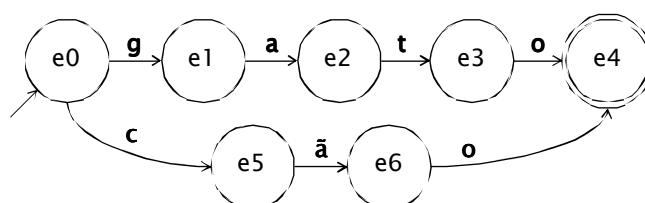
### ► [Pp]reço



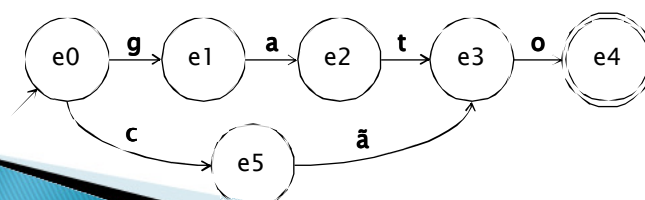
26

## Exemplo

► (gato)|(cão)



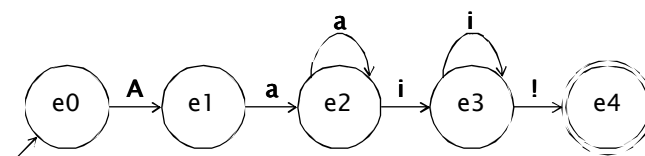
ou



27

## Exemplo

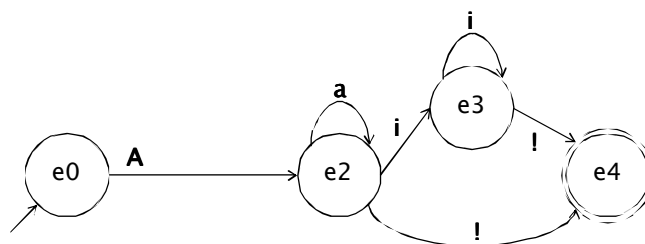
► Aa+i+!



28

## Exemplo

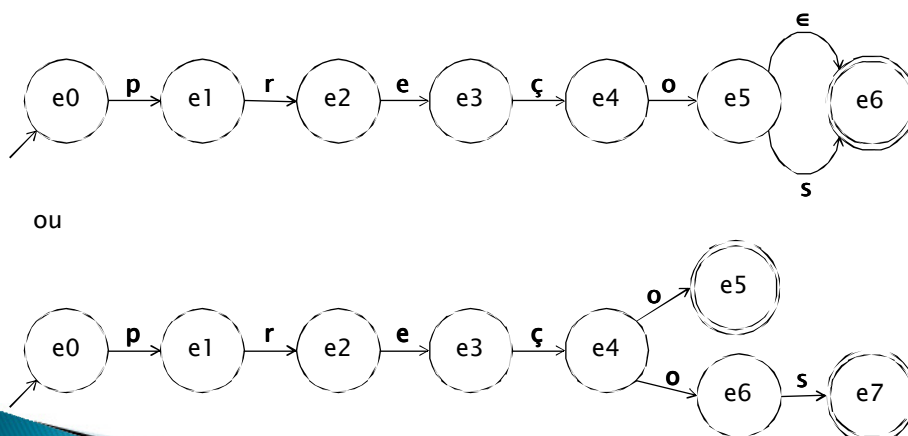
▶  $Aa^*i^*$



29

## Exemplo

▶ preços?



30

## Exercício em duplas

- ▶ 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

- 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

34

## 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

35

## 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

36

## 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

37

## 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 o 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

38


## 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

39

## Exemplo do UNITEX-PB



```

DELAF_PB.dic - Bloco de notas
Arquivo  Editar  Formatar  Exibir  Ajuda
iguarias, iguaria. N:fp
iguaria, iguaria. N:fs
ilaqueações, ilaqueação. N:fp
ilaqueação, ilaqueação. N:fs
ilações, ilação. N:fp
ilação, ilação. N:fs
ilegalidades, ilegalidade. N:fp
ilegalidade, ilegalidade. N:fs
ilegitimidades, ilegitimidade. N:fp
ilegitimidade, ilegitimidade. N:fs
ilhéus, ilha. N:Dmp
ilhéu, ilha. N:Dms
ilhotas, ilha. N:Dfp
ilhota, ilha. N:Dfs
ilhas, ilha. N:fp
ilha, ilha. N:fs
ilhaís, ilha1. N:mp
ilha1, ilha1. N:ms
  
```

## 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

41

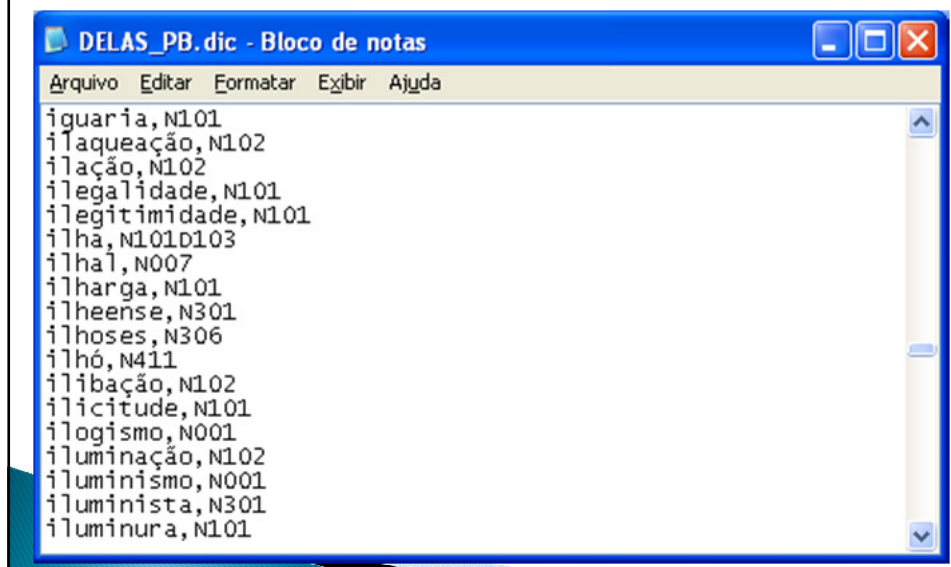
## Análise morfológica

### ▶ 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.

42

## Exemplo do UNITEX–PB



## 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

## Turco – exemplo

(Jurafsky e Martin, 2008)

uygarlaştramadıklarımızdanmışsınızcasma

uygar +laş +nr +ama +dık +lar +ımız +dan +mış +sınız +casma

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

45

## 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

46

## Análise morfológica

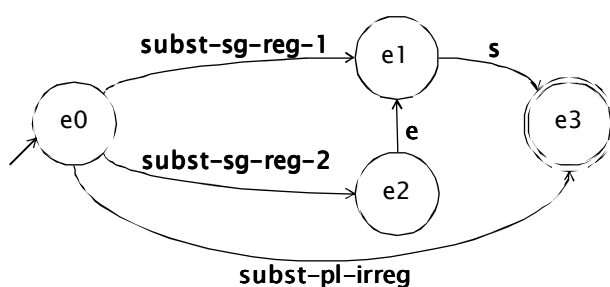
### ► Mapear palavras em seus componentes

- 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

47

## Exemplo simples

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

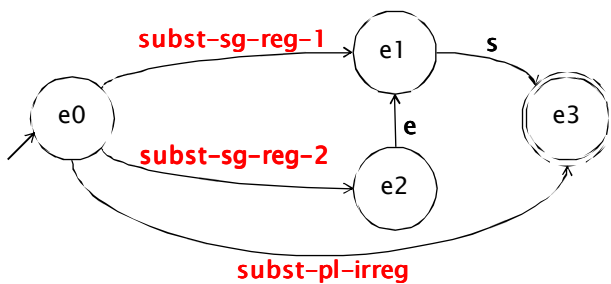


subst-sg-reg-1	subst-sg-reg-2	subst-pl-irreg
casa	flor	lápiz
porta	lar	cópus
...	...	...



## Exemplo simples

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



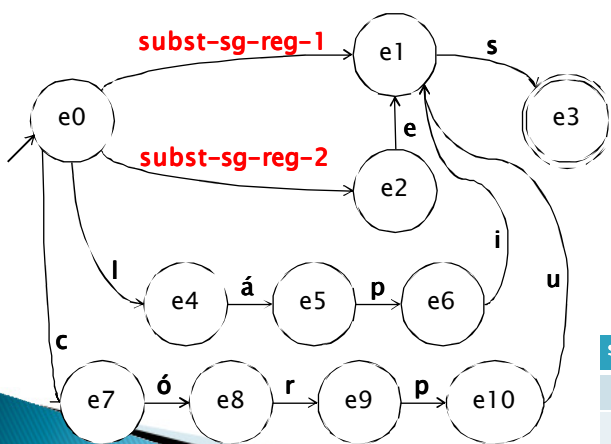
Podem ser substituídos pelos autômatos correspondentes!

Como?

subst-sg-reg-1	subst-sg-reg-2	subst-pl-irreg
casa	flor	lápis
porta	lar	córpus
...	...	...

## Exemplo simples

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



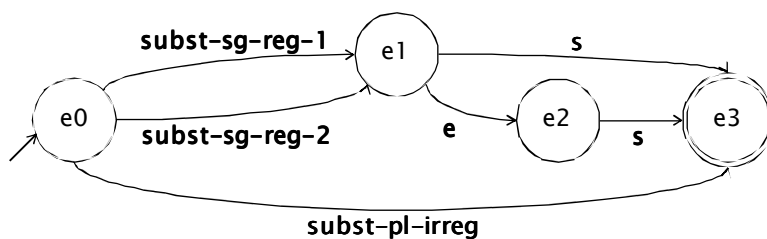
Podem ser substituídos pelos autômatos correspondentes!

Continuem!

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

## Exemplo simples

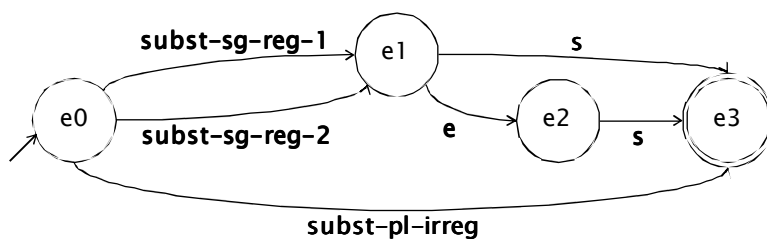
- ▶ Cuidado com **overgeneration!**
  - O que acontece no caso abaixo?



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
...	...	...

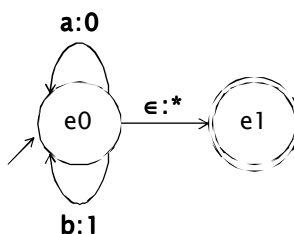
## Análise morfológica

- ▶ 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

53

## Transdutores

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

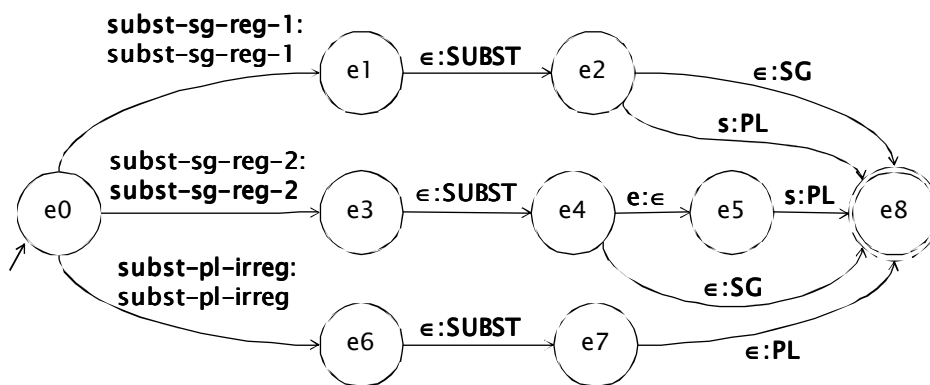


Análise de abba

54

# Transdutores: exemplo

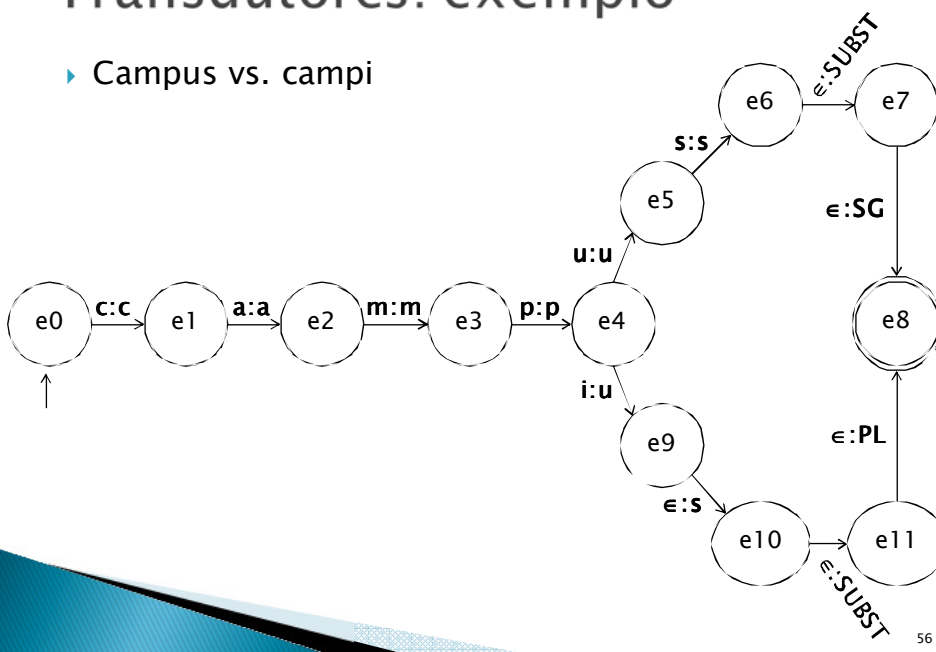
► Releitura do autômato de substantivos



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 em duplas!**
  - Reconhecer número, gênero, raiz e etiqueta morfossintática

57

## 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

58

## 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

59

## 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/funcões similares em seus contextos são de uma mesma classe
  - Algo mais?

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
  - Semântica: substantivos têm uma preferência por objetos, lugares e coisas, adjetivos por propriedades, etc.
  - Pragmática

61

## 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

62

## 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
NPS	Proper noun, plural
PDT	Predeterminer
POS	Possessive ending
PP	Personal pronoun
PP\$	Possessive pronoun
RB	Adverb
RBR	Adverb, comparative
RBS	Adverb, superlative
RP	Particle
SYM	Symbol
TO	to
UH	Interjection
VB	Verb, base form
VBD	Verb, past tense
VBG	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		

64



## 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

65

## 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.

66

## 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

67

## Etiquetagem 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!**

68

## 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.

69

## 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

70

## 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!

71

## Exemplo: EngCG tagger (Voutilainen, 1999)

→ Análise morfológica da sentença (tag correta em negrito)

*Pavlov had shown that salivation.*

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



72

## 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*

73

## 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)

74

## 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

<b>ELAG</b>
16 regras para Adjetivos
30 regras para Advérbios
22 regras para Artigos
12 regras para Substantivos

75

## 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)

76

## Tagging: probabilidades

- ▶ Abordagem antiga, desde a década de 60
- ▶ **Modelo de Markov Oculto** (HMM – *Hidden Markov Model*), um dos mais utilizados
  - Um tipo de inferência bayesiana
  - Tarefa de classificação: dadas algumas observações, quais as classes mais prováveis
    - Tagging: dada uma **sequência de palavras**, qual a **sequência de tags** mais provável

77

## Tagging: probabilidades

- ▶ Tagging: dada uma sequência de palavras, qual a sequência de tags mais provável
- ▶ Queremos a sequência de tags **SeqTags** que **maximize** a probabilidade  $P(\text{SeqTags}|\text{SeqPalavras})$

$$\hat{\text{SeqTags}} = \underset{\text{SeqTags}}{\text{argmax}} P(\text{SeqTags} | \text{SeqPalavras})$$

- ▶ Exemplo: *O menino prefere brincar do que estudar*
  - ▶ SeqTags 1: art subst verbo verbo contração pro verbo
  - ▶ SeqTags 2: art subst verbo verbo contração adv verbo
  - ▶ SeqTags 3: art subst verbo verbo contração conj verbo
  - ▶ SeqTags 4: pro subst verbo verbo contração pro verbo
  - ▶ SeqTags 5: pro subst verbo verbo contração adv verbo
  - ▶ SeqTags 6: pro subst verbo verbo contração conj verbo
- ▶ Qual a melhor sequência de tags, ou seja, qual destas sequências maximiza a probabilidade  $P(\text{SeqTags}|\text{SeqPalavras})$ ?

78

## Tagging: probabilidades

- ▶ Qual a maior probabilidade?

$$\hat{\text{SeqTags}} = \underset{\text{SeqTags}}{\text{argmax}} P(\text{SeqTags} | \text{SeqPalavras})$$

- ▶ P(art subst verbo verbo contração pro verbo | O menino prefere brincar do que estudar)
- ▶ P(art subst verbo verbo contração adv verbo | O menino prefere brincar do que estudar)
- ▶ P(art subst verbo verbo contração conj verbo | O menino prefere brincar do que estudar)
- ▶ P(pro subst verbo verbo contração pro verbo | O menino prefere brincar do que estudar)
- ▶ P(pro subst verbo verbo contração adv verbo | O menino prefere brincar do que estudar)
- ▶ P(pro subst verbo verbo contração conj verbo | O menino prefere brincar do que estudar)

79

## Tagging: probabilidades

- ▶ Como calcular essas probabilidades?

$$P(\text{SeqTags} | \text{SeqPalavras})$$

80



## Tagging: probabilidades

- ▶ Como calcular essas probabilidades?

$$P(\text{SeqTags} | \text{SeqPalavras}) = \frac{P(\text{SeqPalavras} | \text{SeqTags}) \times P(\text{SeqTags})}{P(\text{SeqPalavras})}$$

81

## Tagging: probabilidades

- ▶ Como calcular essas probabilidades?

$$P(\text{SeqTags} | \text{SeqPalavras}) = \frac{P(\text{SeqPalavras} | \text{SeqTags}) \times P(\text{SeqTags})}{P(\text{SeqPalavras})} \longrightarrow \text{constante}$$

$$P(\text{SeqTags} | \text{SeqPalavras}) = P(\text{SeqPalavras} | \text{SeqTags}) \times P(\text{SeqTags})$$

82

## Tagging: probabilidades

- ▶ Como calcular essas probabilidades?

$$P(\text{SeqTags} | \text{SeqPalavras})$$

$$P(\text{SeqTags} | \text{SeqPalavras}) = \frac{P(\text{SeqPalavras} | \text{SeqTags}) \times P(\text{SeqTags})}{P(\text{SeqPalavras})} \longrightarrow \text{constante}$$

$$P(\text{SeqTags} | \text{SeqPalavras}) = P(\text{SeqPalavras} | \text{SeqTags}) \times P(\text{SeqTags})$$

- ▶ Simplificações para facilitar o cálculo
  - ▶ Cada palavra depende apenas de sua tag
  - ▶ Uma tag depende apenas da tag anterior na sentença

83

## Tagging: probabilidades

- ▶ Como calcular essas probabilidades?

$$P(\text{SeqTags} | \text{SeqPalavras})$$

$$P(\text{SeqTags} | \text{SeqPalavras}) = \frac{P(\text{SeqPalavras} | \text{SeqTags}) \times P(\text{SeqTags})}{P(\text{SeqPalavras})} \longrightarrow \text{constante}$$

$$P(\text{SeqTags} | \text{SeqPalavras}) = P(\text{SeqPalavras} | \text{SeqTags}) \times P(\text{SeqTags})$$

- ▶ Simplificações para facilitar o cálculo
  - ▶ Cada palavra depende apenas de sua tag
  - ▶ Uma tag depende apenas da tag anterior na sentença

$$P(\text{SeqTags} | \text{SeqPalavras}) = \prod_{i=1}^{\text{número de palavras}} P(\text{palavra}_i | \text{tag}_i) \times P(\text{tag}_i | \text{tag}_{i-1})$$

84

## Tagging: probabilidades

- ▶ Como calcular essas probabilidades?

$$P(\text{SeqTags} | \text{SeqPalavras})$$

$$P(\text{SeqTags} | \text{SeqPalavras}) = \frac{P(\text{SeqPalavras} | \text{SeqTags}) \times P(\text{SeqTags})}{P(\text{SeqPalavras})} \longrightarrow \text{constante}$$

$$P(\text{SeqTags} | \text{SeqPalavras}) = P(\text{SeqPalavras} | \text{SeqTags}) \times P(\text{SeqTags})$$

- ▶ Simplificações para facilitar o cálculo
  - ▶ Cada palavra depende apenas de sua tag
  - ▶ Uma tag depende apenas da tag anterior na sentença

$$P(\text{SeqTags} | \text{SeqPalavras}) = \prod_{i=1}^{\text{número de palavras}} P(\text{palavra}_i | \text{tag}_i) \times P(\text{tag}_i | \text{tag}_{i-1})$$

Como calcular essas 2 probabilidades?

85

## Tagging: probabilidades

- ▶ Exemplo
  - ▶ Supondo que se usa o Brown Corpus

$$P(\text{SeqTags} | \text{SeqPalavras}) = \prod_{i=1}^{\text{número de palavras}} P(\text{palavra}_i | \text{tag}_i) \times P(\text{tag}_i | \text{tag}_{i-1})$$

$$P(\text{palavra}_i | \text{tag}_i) = P(\text{"is"} | \text{VBZ})$$

tag=VBZ → 21.627 ocorrências no corpus  
 palavra="is" com tag=VBZ → 10.073 ocorrências no corpus

$P(\text{"is"} | \text{VBZ}) = \text{número de vezes de "is" com VBZ} / \text{número de VBZ}$

$P(\text{"is"} | \text{VBZ}) = 10.073 / 21.627 = 0.47$  ou 47%

86

## Tagging: probabilidades

### ▶ Exemplo

- ▶ Supondo que se usa o Brown Corpus

$$P(\text{SeqTags} | \text{SeqPalavras}) = \prod_{i=1}^{\text{número de palavras}} P(\text{palavra}_i | \text{tag}_i) \times P(\text{tag}_i | \text{tag}_{i-1})$$

$$P(\text{tag}_i | \text{tag}_{i-1}) = P(\text{NN} | \text{DT})$$

tag<sub>i-1</sub>=DT → 116.454 ocorrências no corpus  
tag<sub>i</sub>=NN com tag<sub>i-1</sub>=DT → 56.509 ocorrências no corpus

P(NN|DT) = número de vezes de NN precedido por DT / número de DT  
P(NN|DT) = 56.509 / 116.454 = 0.49 ou 49%

87

## Tagging: probabilidades

### ▶ Exemplo

- ▶ P(art subst verbo verbo contração pro verbo | O menino prefere brincar do que estudar)

$$P(\text{SeqTags} | \text{SeqPalavras}) = \prod_{i=1}^{\text{número de palavras}} P(\text{palavra}_i | \text{tag}_i) \times P(\text{tag}_i | \text{tag}_{i-1})$$

- ▶ = P(O|art) x P(art|<início da sentença>) x  
P(menino|subst) x P(subst|art) x  
P(prefere|verbo) x P(verbo|subst) x  
P(brincar|verbo) x P(verbo|verbo) x  
P(do|contração) x P(contração|verbo) x  
P(que|pro) x P(pro|contração) x  
P(estudar|verbo) x P(verbo|pro)

- ▶ Faz-se isso para **todas as possíveis sequências de tags** (com probabilidades aprendidas de corpus)

- ▶ **A sequência com maior probabilidade é escolhida**

88

## Tagging: probabilidades

- ▶ Considerando a probabilidade abaixo
  - ▶  $P(\text{art subst verbo verbo contração pro verbo} \mid \text{O menino prefere brincar do que estudar}) =$   
 $P(\text{O}|\text{art}) \times P(\text{art}|\langle \text{início da sentença} \rangle) \times$   
 $P(\text{menino}|\text{subst}) \times P(\text{subst}|\text{art}) \times$   
 $P(\text{prefere}|\text{verbo}) \times P(\text{verbo}|\text{subst}) \times$   
 $P(\text{brincar}|\text{verbo}) \times P(\text{verbo}|\text{verbo}) \times$   
 $P(\text{do}|\text{contração}) \times P(\text{contração}|\text{verbo}) \times$   
 $P(\text{que}|\text{pro}) \times P(\text{pro}|\text{contração}) \times$   
 $P(\text{estudar}|\text{verbo}) \times P(\text{verbo}|\text{pro})$
  - ▶ Por que ela é **provavelmente menor** do que a abaixo – que é a interpretação correta?
    - ▶  $P(\text{art subst verbo verbo contração conj verbo} \mid \text{O menino prefere brincar do que estudar}) =$   
 $P(\text{O}|\text{art}) \times P(\text{art}|\langle \text{início da sentença} \rangle) \times$   
 $P(\text{menino}|\text{subst}) \times P(\text{subst}|\text{art}) \times$   
 $P(\text{prefere}|\text{verbo}) \times P(\text{verbo}|\text{subst}) \times$   
 $P(\text{brincar}|\text{verbo}) \times P(\text{verbo}|\text{verbo}) \times$   
 $P(\text{do}|\text{contração}) \times P(\text{contração}|\text{verbo}) \times$   
 $P(\text{que}|\text{conj}) \times P(\text{conj}|\text{contração}) \times$   
 $P(\text{estudar}|\text{verbo}) \times P(\text{verbo}|\text{conj})$

89

## Tagging: probabilidades

- ▶ Considerando a probabilidade abaixo
  - ▶  $P(\text{art subst verbo verbo contração pro verbo} \mid \text{O menino prefere brincar do que estudar}) =$   
 $P(\text{O}|\text{art}) \times P(\text{art}|\langle \text{início da sentença} \rangle) \times$   
 $P(\text{menino}|\text{subst}) \times P(\text{subst}|\text{art}) \times$   
 $P(\text{prefere}|\text{verbo}) \times P(\text{verbo}|\text{subst}) \times$   
 $P(\text{brincar}|\text{verbo}) \times P(\text{verbo}|\text{verbo}) \times$   
 $P(\text{do}|\text{contração}) \times P(\text{contração}|\text{verbo}) \times$   
 $P(\text{que}|\text{pro}) \times P(\text{pro}|\text{contração}) \times$   
 $P(\text{estudar}|\text{verbo}) \times P(\text{verbo}|\text{pro})$
  - ▶ Por que ela é **provavelmente menor** do que a abaixo – que é a interpretação correta?
    - ▶  $P(\text{art subst verbo verbo contração conj verbo} \mid \text{O menino prefere brincar do que estudar}) =$   
 $P(\text{O}|\text{art}) \times P(\text{art}|\langle \text{início da sentença} \rangle) \times$   
 $P(\text{menino}|\text{subst}) \times P(\text{subst}|\text{art}) \times$   
 $P(\text{prefere}|\text{verbo}) \times P(\text{verbo}|\text{subst}) \times$   
 $P(\text{brincar}|\text{verbo}) \times P(\text{verbo}|\text{verbo}) \times$   
 $P(\text{do}|\text{contração}) \times P(\text{contração}|\text{verbo}) \times$   
 $P(\text{que}|\text{conj}) \times P(\text{conj}|\text{contração}) \times$   
 $P(\text{estudar}|\text{verbo}) \times P(\text{verbo}|\text{conj})$

Esses termos devem ser mais prováveis do que os correspondentes na interpretação errada

90

## Tagging: probabilidades

### ▶ Modelo de Markov oculto

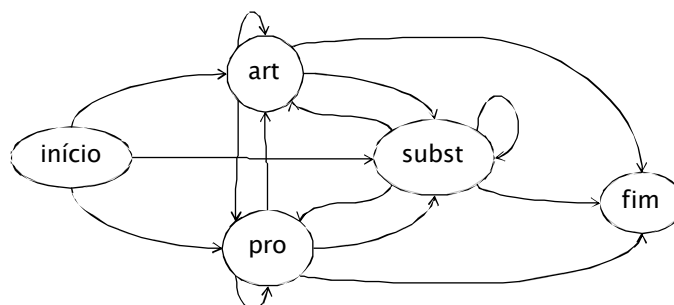
- Modelam-se eventos observados (palavras) e eventos não observados, ou seja, ocultos (tags)
- Como dito anteriormente, tipo especial de autômato
  - Probabilidades nos arcos (transições):  $P(\text{tag}_i | \text{tag}_{i-1})$
  - Probabilidades nos nós:  $P(\text{palavra}_i | \text{tag}_i)$

91

## Tagging: probabilidades

### ▶ Modelo de Markov oculto

- Exemplo hipotético



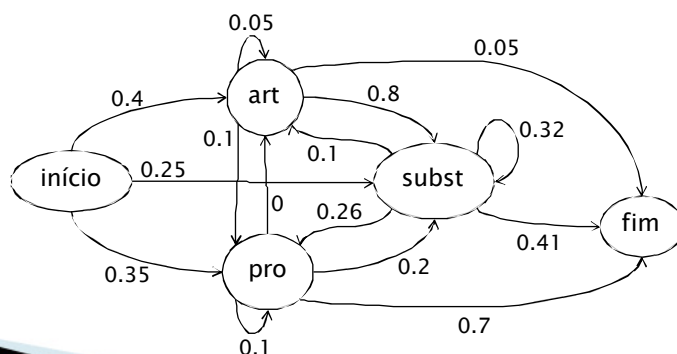
92

## Tagging: probabilidades

### ► Modelo de Markov oculto

#### ◦ Exemplo hipotético

A soma das probabilidades dos arcos que saem de cada nó devem somar 1



93

## Tagging: probabilidades

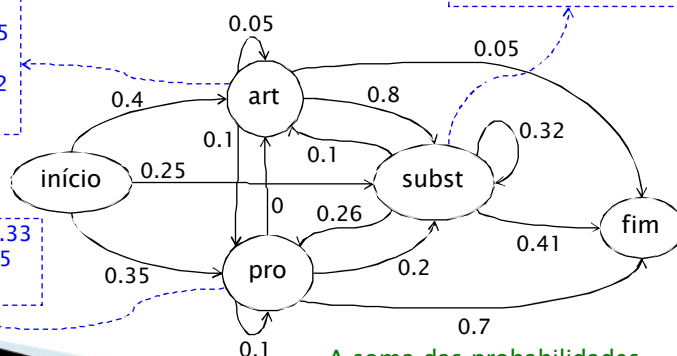
### ► Modelo de Markov oculto

#### ◦ Exemplo hipotético

$P(o|art)=0.2$   
 $P(os|art)=0.25$   
 $P(a|art)=0.2$   
 $P(as|art)=0.12$   
 ...

$P(menino|subst)=0.1$   
 $P(casa|subst)=0.04$   
 ...

$P(que|pro)=0.33$   
 $P(ele|pro)=0.5$   
 ...



A soma das probabilidades associadas a cada nó devem somar 1

## Tagging: probabilidades

- ▶ **Taggers atuais** usam normalmente **2 tags anteriores** como contexto
  - $P(\text{tag}_i | \text{tag}_{i-1}, \text{tag}_{i-2})$
- ▶ É preciso **otimizar a busca por sequências de tags** mais prováveis, senão pode haver explosão combinatória
  - **Programação dinâmica** (estudaremos logo)
- ▶ É preciso **lidar com probabilidades muito baixas ou próximas de zero**
  - Por exemplo, situações em que  $P(\text{tag}_i | \text{tag}_{i-1}, \text{tag}_{i-2}) \approx 0$
  - Solução usual: combinar probabilidades ponderadas
    - $P(\text{tag}_i | \text{tag}_{i-1}, \text{tag}_{i-2}) = p_1 * P(\text{tag}_i | \text{tag}_{i-1}, \text{tag}_{i-2}) + p_2 * P(\text{tag}_i | \text{tag}_{i-1}) + p_3 * P(\text{tag}_i)$ , com  $p_1 + p_2 + p_3 = 1$

95

## Tagging: Transformation-Based

- ▶ *Transformation-Based tagging*
  - Aplicação de TBL (*Transformation-Based Learning*), da linha de aprendizado de máquina
  - **Regras são aprendidas e aprimoradas** automaticamente
    - Várias iterações
    - A cada iteração, o processo melhora
    - Ao estabilizar, fim do processo de aprendizado

96



## Tagging: Transformation-Based

### ▶ Transformation-Based tagging

- **Processo básico** com base em um **cópus anotado manualmente**
  - Inicialmente, **anota-se um cópus automaticamente**, assumindo-se que a tag de uma palavra é a sua **tag mais frequente** (segundo um cópus/léxico)
  - Verificam-se os **erros cometidos** (comparando-se com a anotação humana correspondente) e, dentre todas as possibilidades de correção, **monta-se uma regra de correção** com maior precisão
  - Aplica-se essa **regra nova em todo o cópus**
  - Verificam-se novamente os **erros cometidos** e monta-se uma segunda **regra de correção** com maior precisão
  - E assim por diante, **até não se obter mais melhora de performance**

97

## Tagging: Transformation-Based

### ▶ Exemplo

#### 1. Inicialmente, anotação com base em frequência

... is/VBZ expected/VBN to/TO **race/NN** tomorrow/NN  
 ... the/DT race/NN for/IN outer/JJ space/NN  
 Book/VB the/DT **flight/VB** to/TO...

#### 2. Verificando-se os erros, aprende-se uma nova regra

**Troque NN para VB quando a tag anterior é TO**

#### 3. Corrige-se a etiquetagem anterior

... is/VBZ expected/VBN to/TO **race/VB** tomorrow/NN  
 ... the/DT race/NN for/IN outer/JJ space/NN  
 Book/VB the/DT **flight/VB** to/TO...

98

## Tagging: Transformation-Based

### ▶ Exemplo

4. Aprende-se uma nova regra com base nos erros existentes

Troque VB para NN quando a tag anterior é DT e a posterior é TO

5. Corrige-se a etiquetação anterior

... is/VBZ expected/VBN to/TO **race/VB** tomorrow/NN

... the/DT race/NN for/IN outer/JJ space/NN

Book/VB the/DT **flight/NN** to/TO...

6. E assim por diante

99

## Tagging: Transformation-Based

### ▶ Exemplo

- **Resultado** do processo

Regra 1: etiqüete as palavras com suas tags mais frequentes

Regra 2: troque NN para VB quando a tag anterior é TO

Regra 3: troque VB para NN quando a tag anterior é DT e a posterior é TO

...

100

## Tagging: Transformation-Based

### ▶ Transformation-Based tagging

- Ao final do processo, há um **conjunto de regras ordenadas** que devem ser aplicadas sequencialmente para etiquetar um novo texto
- Como há muitas regras possíveis de serem aprendidas, costuma-se **limitar as estruturas aceitas de regras**
  - Troque a tag da palavra corrente de A para B se a palavra anterior tem a tag X
  - Troque a tag da palavra corrente de A para B se a palavra anterior tem a tag X e a posterior tem a tag Y
  - Etc.

Caso contrário, o que acontece?

101

## Humanos e máquinas

### ▶ Palavras/morfologia e léxico mental

- Nem a listagem exaustiva, nem todas as regras de flexão/derivação
- Há indícios de que humanos armazenam em seu léxico mental os lemas das palavras e também algumas formas plenas
  - Stanners et al. (1979): são mantidas separadamente as palavras *happy* e *happiness*, mas somente o verbo *pour*, sem suas flexões

### ▶ Morfossintaxe

- Experimentos mostram que humanos discordam em 3-4% das tags
  - Melhores taggers com 97%, normalmente com problemas justamente onde os humanos discordam
- Voutilainen (1995) mostra que humanos atingem 100% se se permite que discutam as tags com problemas

102