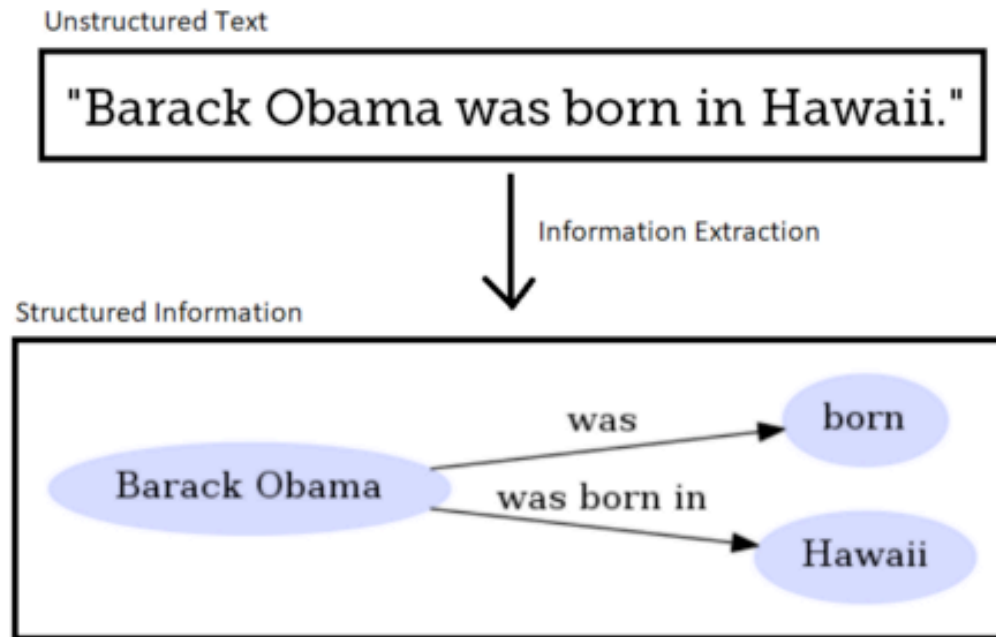


Information Extraction

Luciano Barbosa

Objetivo

- Extrair estrutura a partir de dados não estruturados



Sequence Labeling

- Objetivo: atribuir um dado rótulo a cada palavra de um sentença
- Rótulos dependem de outras palavras das sequência (não é i.i.d)

Part Of Speech Tagging

- Objetivo: atribuir a classe gramatical a cada palavra de uma sentença (substantivo, adjetivo, verbo etc)
- Útil para parsing sintático e desambiguação de palavras
John saw the saw and decided to take it to the table.
PN V Det N Con V Part V Pro Prep Det N

Sequence Tagging

- Objetivo: atribuir rótulos a palavras de uma sentença (pessoa, local, empresa)
- Named entity recognition: identifica nomes de pessoas, locais no texto

people organizations places

- Michael Dell is the CEO of Dell Computer Corporation and lives in Austin Texas.

- Extrair partes de informação relevante para uma dada aplicação

make model year mileage price

- For sale, 2002 Toyota Prius, 20,000 mi, \$15K or best offer. Available starting July 30, 2006.

Semantic Role Labeling

- Determina o papel semântico de cada noun phrase que é argumento do verbo

agent patient source destination instrument

- John drove Mary from Austin to Dallas in his Toyota Prius.
- The hammer broke the window.

Bioinformática

- Rotular sequências genéticas

extron intron

— AGCTAACGTTCGATACGGATTACAGCCT

Principais POS (Inglês)

	Tag	Description	Example
Open Class	ADJ	Adjective: noun modifiers describing properties	<i>red, young, awesome</i>
	ADV	Adverb: verb modifiers of time, place, manner	<i>very, slowly, home, yesterday</i>
	NOUN	words for persons, places, things, etc.	<i>algorithm, cat, mango, beauty</i>
	VERB	words for actions and processes	<i>draw, provide, go</i>
	PROPN	Proper noun: name of a person, organization, place, etc..	<i>Regina, IBM, Colorado</i>
	INTJ	Interjection: exclamation, greeting, yes/no response, etc.	<i>oh, um, yes, hello</i>
Closed Class Words	ADP	Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation	<i>in, on, by under</i>
	AUX	Auxiliary: helping verb marking tense, aspect, mood, etc.,	<i>can, may, should, are</i>
	CCONJ	Coordinating Conjunction: joins two phrases/clauses	<i>and, or, but</i>
	DET	Determiner: marks noun phrase properties	<i>a, an, the, this</i>
	NUM	Numeral	<i>one, two, first, second</i>
	PART	Particle: a preposition-like form used together with a verb	<i>up, down, on, off, in, out, at, by</i>
	PRON	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others</i>
Other	SCONJ	Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement	<i>that, which</i>
	PUNCT	Punctuation	<i>; , ()</i>
	SYM	Symbols like \$ or emoji	<i>\$, %</i>
	X	Other	<i>asdf, qwfg</i>

Tipos de POS

- Closed class:
 - Preposições e pronomes
 - Tendem a ser curtos
 - Alta frequência
- Open class:
 - Substantivos, verbos, adjetivos e advérbios
 - Constantemente sendo criados

POS no Penn Treebank

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coord. conj.	<i>and, but, or</i>	NNP	proper noun, sing.	<i>IBM</i>	TO	“to”	<i>to</i>
CD	cardinal number	<i>one, two</i>	NNPS	proper noun, plu.	<i>Carolinas</i>	UH	interjection	<i>ah, oops</i>
DT	determiner	<i>a, the</i>	NNS	noun, plural	<i>llamas</i>	VB	verb base	<i>eat</i>
EX	existential ‘there’	<i>there</i>	PDT	predeterminer	<i>all, both</i>	VBD	verb past tense	<i>ate</i>
FW	foreign word	<i>mea culpa</i>	POS	possessive ending	<i>’s</i>	VBG	verb gerund	<i>eating</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	PRP	personal pronoun	<i>I, you, he</i>	VBN	verb past partici- ple	<i>eaten</i>
JJ	adjective	<i>yellow</i>	PRP\$	possess. pronoun	<i>your, one’s</i>	VBP	verb non-3sg-pr	<i>eat</i>
JJR	comparative adj	<i>bigger</i>	RB	adverb	<i>quickly</i>	VBZ	verb 3sg pres	<i>eats</i>
JJS	superlative adj	<i>wildest</i>	RBR	comparative adv	<i>faster</i>	WDT	wh-determ.	<i>which, that</i>
LS	list item marker	<i>1, 2, One</i>	RBS	superlatv. adv	<i>fastest</i>	WP	wh-pronoun	<i>what, who</i>
MD	modal	<i>can, should</i>	RP	particle	<i>up, off</i>	WP\$	wh-possess.	<i>whose</i>
NN	sing or mass noun	<i>llama</i>	SYM	symbol	<i>+, %, &</i>	WRB	wh-adverb	<i>how, where</i>

Exemplos

There/**PRO/EX** are/**VERB/VBP** 70/**NUM/CD** children/**NOUN/NNS**
there/**ADV/RB** ./**PUNC/**.

Preliminary/**ADJ/JJ** findings/**NOUN/NNS** were/**AUX/VBD** reported/**VERB/VBN**
in/**ADP/IN** today/**NOUN/NN** 's/**PART/POS** New/**PROPN/NNP**
England/**PROPN/NNP** Journal/**PROPN/NNP** of/**ADP/IN** Medicine/**PROPN/NNP**

Part-of-Speech Tagging

- Atribuir a classe gramatical a cada palavra de uma sentença (substantivo, adjetivo, verbo etc)
- Tarefa de desambiguação: palavras podem ter mais de uma class gramatical
 - Ex: book, that etc
- Classe mais frequente: 92% de acurácia
- Estado da arte: 97%

Types:		WSJ	Brown
Unambiguous	(1 tag)	44,432 (86%)	45,799 (85%)
Ambiguous	(2+ tags)	7,025 (14%)	8,050 (15%)
Tokens:			
Unambiguous	(1 tag)	577,421 (45%)	384,349 (33%)
Ambiguous	(2+ tags)	711,780 (55%)	786,646 (67%)

Named Entity Recognition (Information Extraction)

- Named entity: tudo que se refere a um nome próprio (regra geral)

Type	Tag	Sample Categories	Example sentences
People	PER	people, characters	Turing is a giant of computer science.
Organization	ORG	companies, sports teams	The IPCC warned about the cyclone.
Location	LOC	regions, mountains, seas	Mt. Sanitas is in Sunshine Canyon .
Geo-Political Entity	GPE	countries, states	Palo Alto is raising the fees for parking.

- Pode ser qualquer entidade: produto, doenças etc
- Natural language understanding: Q&A, chatbot
- Dificuldades:
 - Encontrar o pedaço do texto que contém a entidade
 - Ambiguidade: JFK (pessoa ou aeroporto)

BIO Tagging

- Convenção para rotulagem de sequência

Words	IO Label	BIO Label	BIOES Label
Jane	I-PER	B-PER	B-PER
Villanueva	I-PER	I-PER	E-PER
of	O	O	O
United	I-ORG	B-ORG	B-ORG
Airlines	I-ORG	I-ORG	I-ORG
Holding	I-ORG	I-ORG	E-ORG
discussed	O	O	O
the	O	O	O
Chicago	I-LOC	B-LOC	S-LOC
route	O	O	O
.	O	O	O

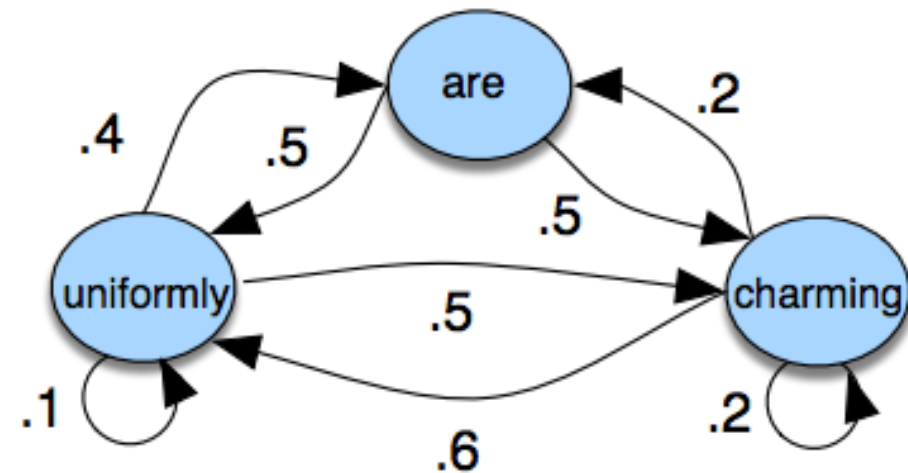
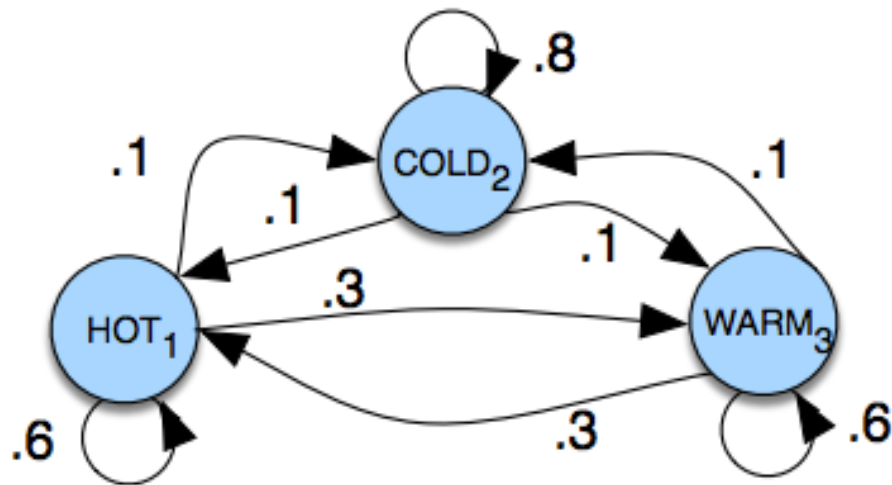
HMM POS Tagging

- Modelo probabilístico sequencial
- Computa a probabilidade para possíveis sequências de rótulos
- Escolhe a melhor sequência

Markov Chain

- Probabilidade do próximo rótulo só depende do anterior

$$P(q_i = a | q_1 \dots q_{i-1}) = P(q_i = a | q_{i-1})$$



Markov Chain: Componentes

$$Q = q_1 q_2 \dots q_N$$

a set of N states

$$A = a_{11} a_{12} \dots a_{N1} \dots a_{NN}$$

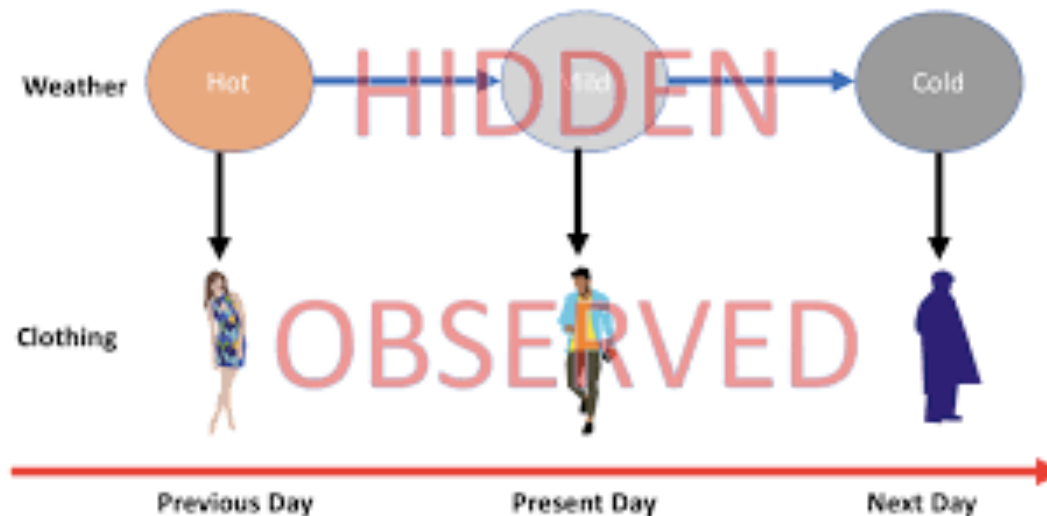
a **transition probability matrix** A , each a_{ij} representing the probability of moving from state i to state j , s.t.
 $\sum_{j=1}^n a_{ij} = 1 \quad \forall i$

$$\pi = \pi_1, \pi_2, \dots, \pi_N$$

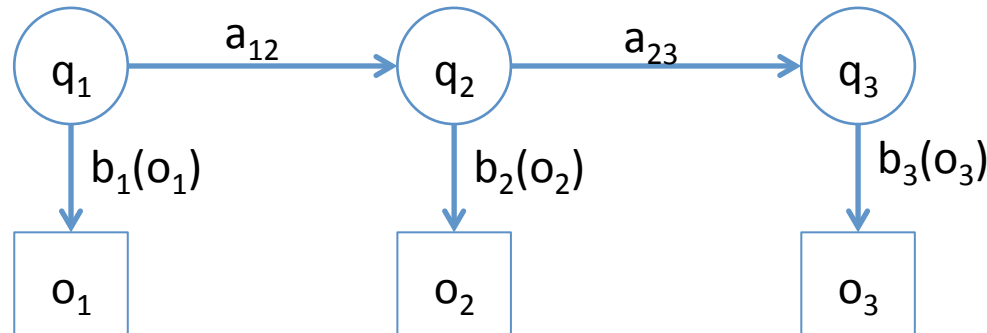
an **initial probability distribution** over states. π_i is the probability that the Markov chain will start in state i . Some states j may have $\pi_j = 0$, meaning that they cannot be initial states. Also, $\sum_{i=1}^n \pi_i = 1$

Hidden Markov Model

- Rótulos estão escondidos (hidden)
- Observa palavras
- Inferir rótulos (ex. POS) da sequência de palavras



Hidden Markov Model: Componentes



$Q = q_1 q_2 \dots q_N$	a set of N states
$A = a_{11} \dots a_{ij} \dots a_{NN}$	a transition probability matrix A , each a_{ij} representing the probability of moving from state i to state j , s.t. $\sum_{j=1}^N a_{ij} = 1 \quad \forall i$
$O = o_1 o_2 \dots o_T$	a sequence of T observations , each one drawn from a vocabulary $V = v_1, v_2, \dots, v_V$
$B = b_i(o_t)$	a sequence of observation likelihoods , also called emission probabilities , each expressing the probability of an observation o_t being generated from a state q_i
$\pi = \pi_1, \pi_2, \dots, \pi_N$	an initial probability distribution over states. π_i is the probability that the Markov chain will start in state i . Some states j may have $\pi_j = 0$, meaning that they cannot be initial states. Also, $\sum_{i=1}^n \pi_i = 1$

Markov Assumption: $P(q_i | q_1, \dots, q_{i-1}) = P(q_i | q_{i-1})$

Output Independence: $P(o_i | q_1, \dots, q_i, \dots, q_T, o_1, \dots, o_i, \dots, o_T) = P(o_i | q_i)$

HMM Tagger

- Matrix A: probabilidades de transição das tags

$$P(t_i|t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$

WSJ corpus

$$P(VB|MD) = \frac{C(MD, VB)}{C(MD)} = \frac{10471}{13124} = .80$$

verbo modal

HMM Tagger

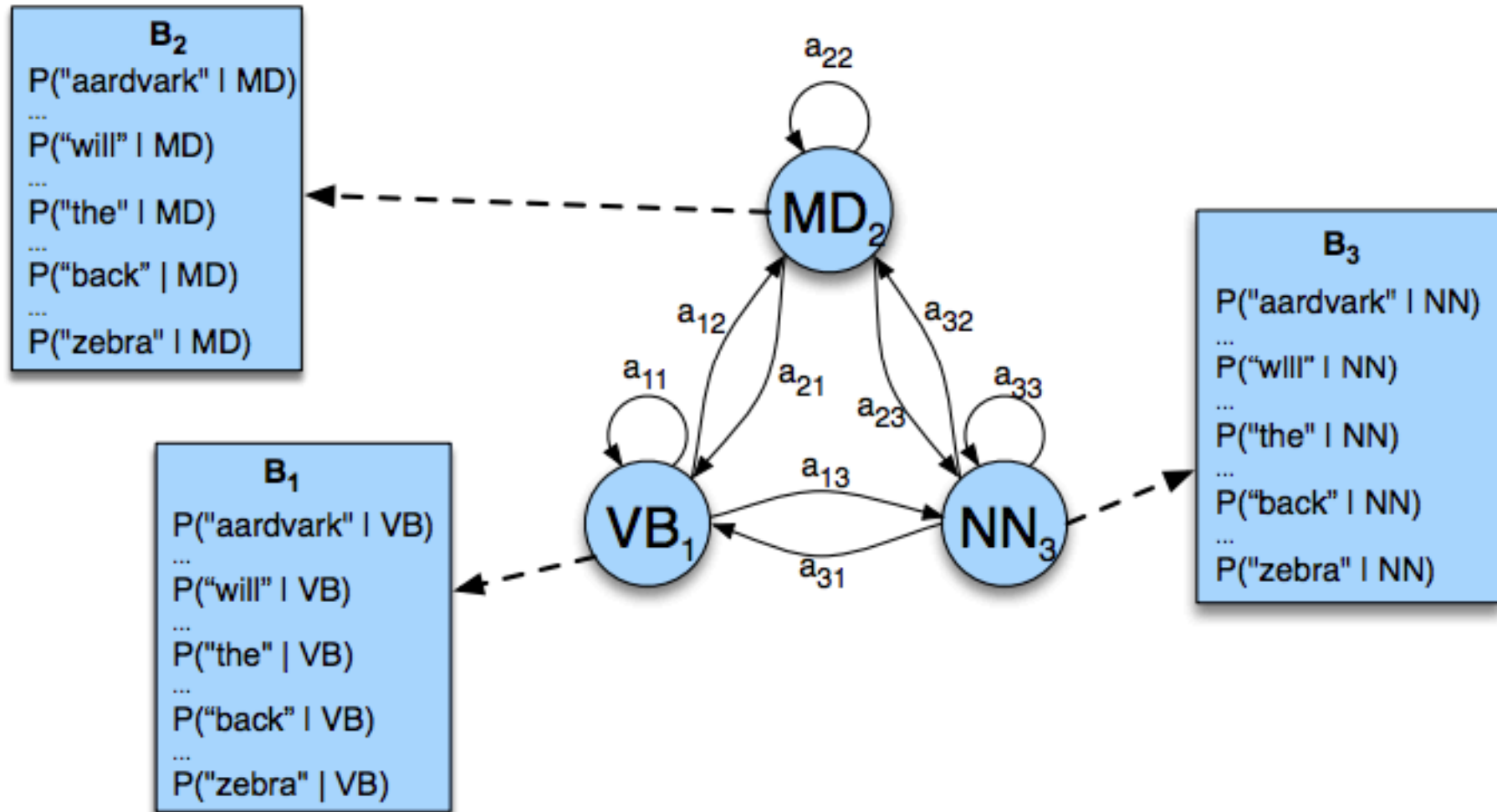
- Matrix B: probabilidades de uma palavra associada a uma tag

$$P(w_i|t_i) = \frac{C(t_i, w_i)}{C(t_i)}$$

WSJ corpus

$$P(will|MD) = \frac{C(MD, will)}{C(MD)} = \frac{4046}{13124} = .31$$

HMM Tagger: Exemplo



Decoding

- Dado um HMM como entrada: matrizes (A,B) e uma sequência de observações (palavras), encontrar a sequência de estados mais provável

$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(t_1 \dots t_n | w_1 \dots w_n)$$

$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} \frac{P(w_1 \dots w_n | t_1 \dots t_n) P(t_1 \dots t_n)}{P(w_1 \dots w_n)}$$

$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(w_1 \dots w_n | t_1 \dots t_n) P(t_1 \dots t_n)$$

$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(t_1 \dots t_n | w_1 \dots w_n) \approx \operatorname{argmax}_{t_1 \dots t_n} \prod_{i=1}^n \overbrace{P(w_i | t_i)}^{\text{emission}} \overbrace{P(t_i | t_{i-1})}^{\text{transition}}$$

Matriz B
Matriz A

Viterbi Algorithm

- Input: Janet will back the bill
- Output: Janet/NNP will/MD back/VB the/DT bill/NN

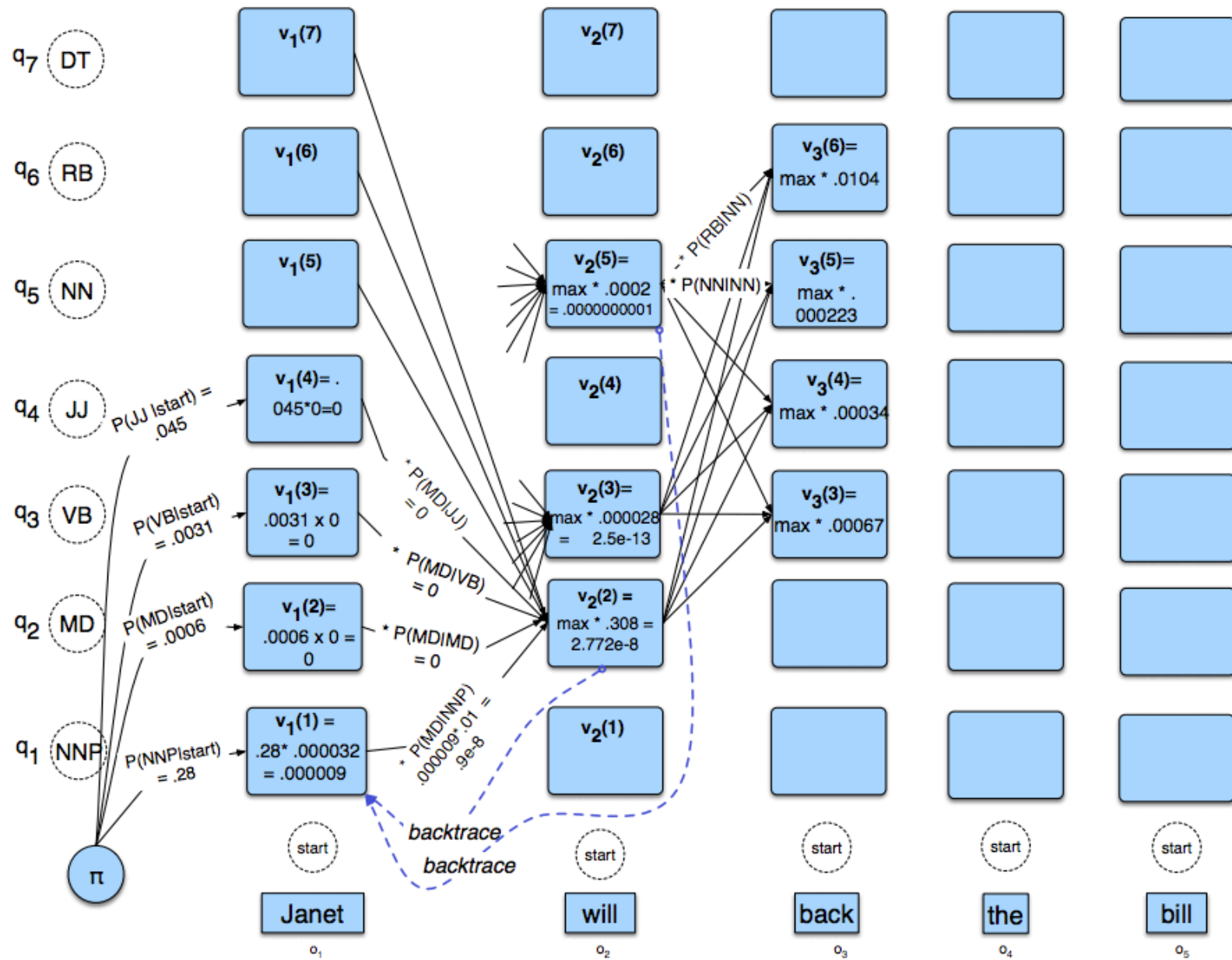
	NNP	MD	VB	JJ	NN	RB	DT
<s>	0.2767	0.0006	0.0031	0.0453	0.0449	0.0510	0.2026
NNP	0.3777	0.0110	0.0009	0.0084	0.0584	0.0090	0.0025
MD	0.0008	0.0002	0.7968	0.0005	0.0008	0.1698	0.0041
VB	0.0322	0.0005	0.0050	0.0837	0.0615	0.0514	0.2231
JJ	0.0366	0.0004	0.0001	0.0733	0.4509	0.0036	0.0036
NN	0.0096	0.0176	0.0014	0.0086	0.1216	0.0177	0.0068
RB	0.0068	0.0102	0.1011	0.1012	0.0120	0.0728	0.0479
DT	0.1147	0.0021	0.0002	0.2157	0.4744	0.0102	0.0017

Matriz A (WSJ corpus)

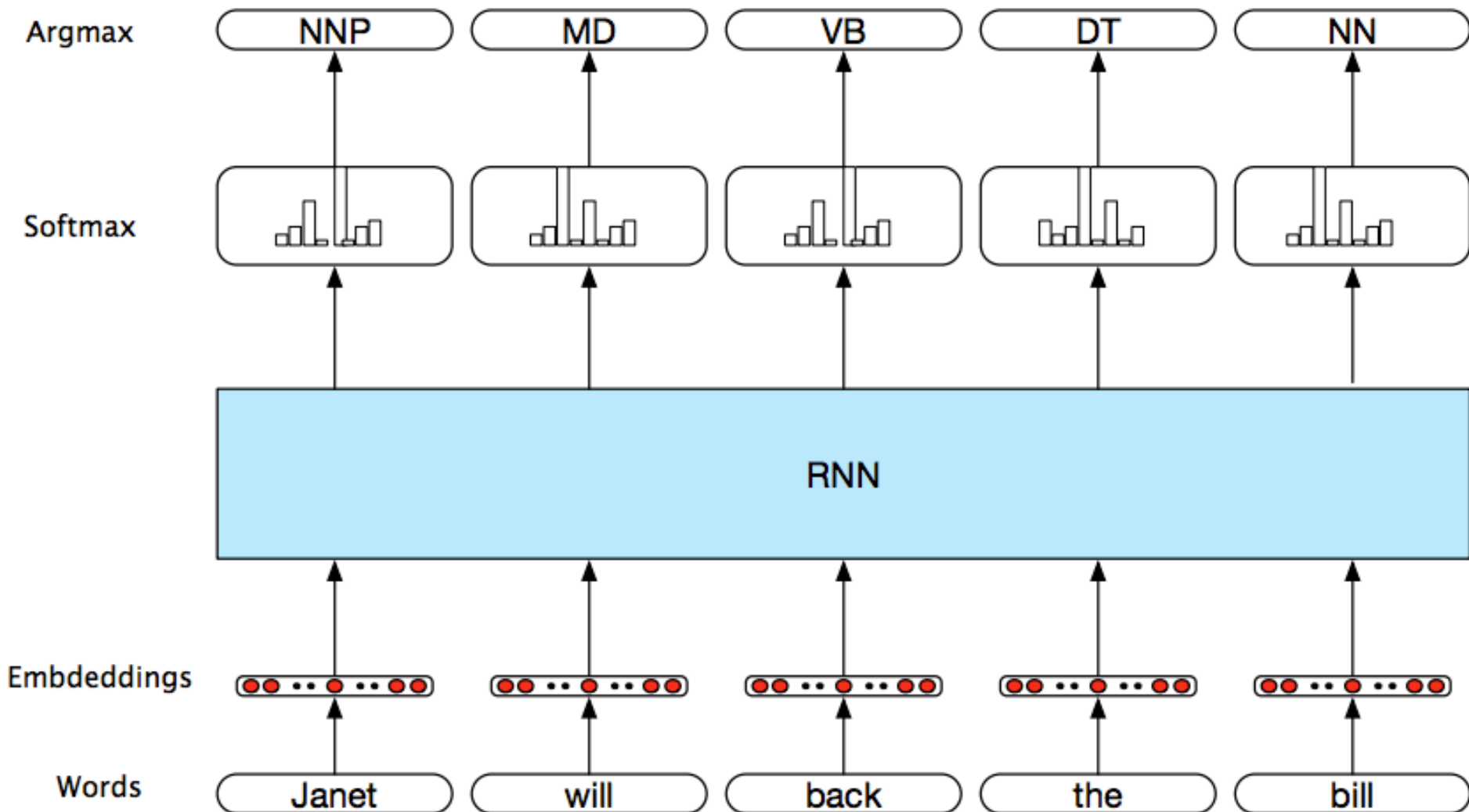
Viterbi Algorithm

	Janet	will	back	the	bill
NNP	0.000032	0	0	0.000048	0
MD	0	0.308431	0	0	0
VB	0	0.000028	0.000672	0	0.000028
JJ	0	0	0.000340	0	0
NN	0	0.000200	0.000223	0	0.002337
RB	0	0	0.010446	0	0
DT	0	0	0	0.506099	0

Matriz B (WSJ corpus)



RNN para POS Tagging

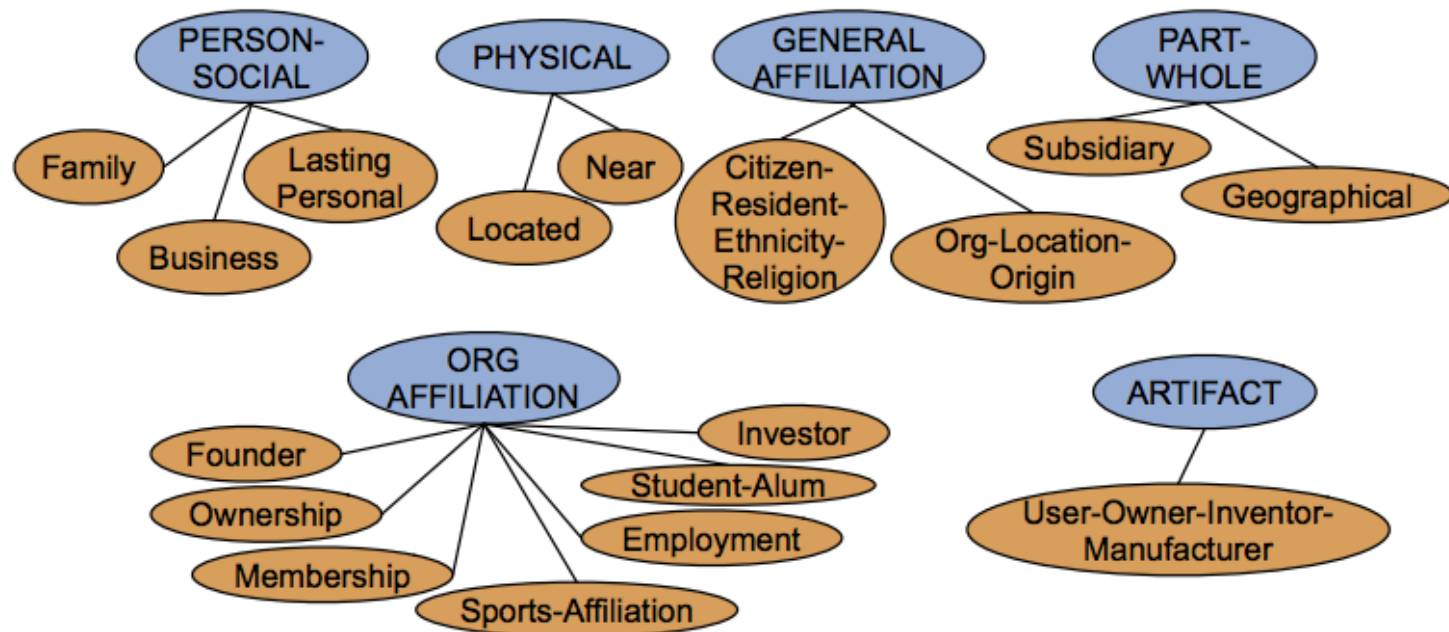


Relation Extraction

- Encontrar e classificar relações entre entidades

Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY \$6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PER Tim Wagner] said. [ORG United], a unit of [ORG UAL Corp.], said the increase took effect [TIME Thursday] and applies to most routes where it competes against discount carriers, such as [LOC Chicago] to [LOC Dallas] and [LOC Denver] to [LOC San Francisco].

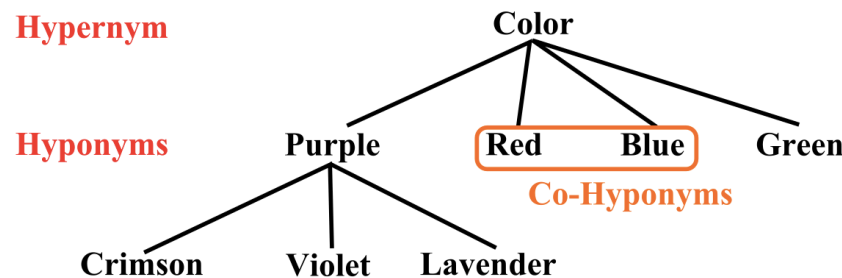
Exemplos de Relações



Relations	Types	Examples
Physical-Located	PER-GPE	He was in Tennessee
Part-Whole-Subsidiary	ORG-ORG	XYZ , the parent company of ABC
Person-Social-Family	PER-PER	Yoko 's husband John
Org-AFF-Founder	PER-ORG	Steve Jobs , co-founder of Apple...

Extração Baseada em Padrões

- Hearst para extração de hyponyms



- Exemplo

Agar is a substance prepared from a mixture of red algae, such as Gelidium, for laboratory or industrial use.

- Padrão:

NP_0 such as $NP_1\{,NP_2\dots,(and|or)NP_i\}, i \geq 1$



hyponym(Gelidium, red algae)

Extração Baseada em Padrões

- Hearst para extração de hyponyms

NP {, NP}* {,} (and or) other NP _H	temples, treasuries, and other important civic buildings
NP _H such as {NP,}* {(or and)} NP	red algae such as Gelidium
such NP _H as {NP,}* {(or and)} NP	such authors as Herrick, Goldsmith, and Shakespeare
NP _H {,} including {NP,}* {(or and)} NP	common-law countries , including Canada and England
NP _H {,} especially {NP,}* {(or and)} NP	European countries , especially France, England, and Spain

Extração Baseada em ML

1. Definem-se as relações e entidades a serem extraídas
2. Anotam-se exemplos para treinamento

function FINDRELATIONS(*words*) **returns** *relations*

relations \leftarrow nil

entities \leftarrow FINDENTITIES(*words*)

forall entity pairs $\langle e1, e2 \rangle$ **in** *entities* **do**

if RELATED?(*e1*, *e2*)

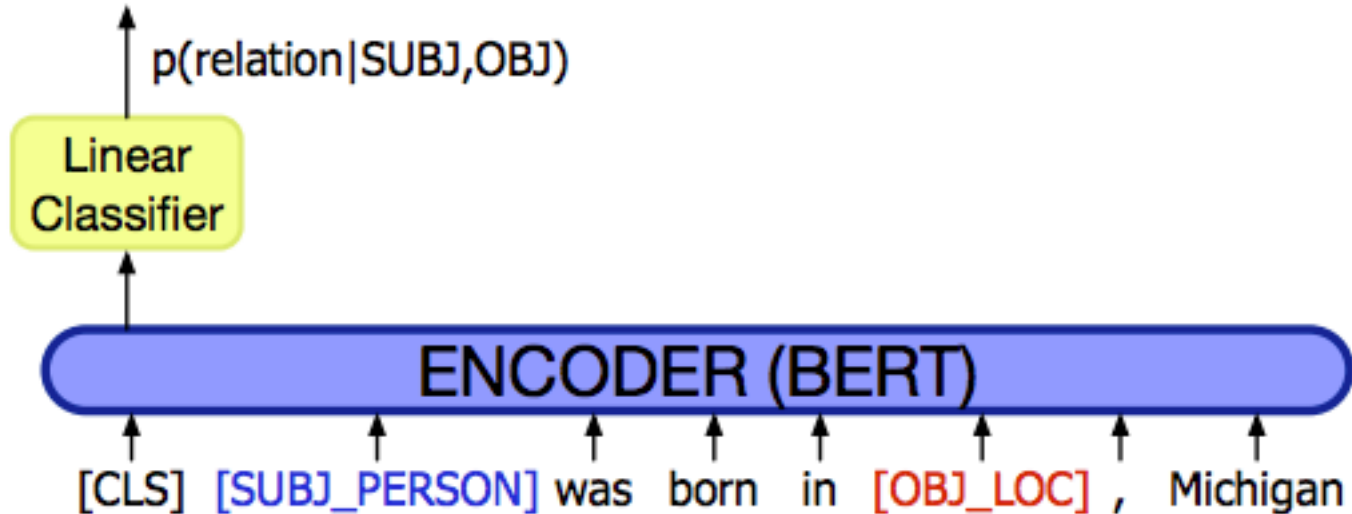
relations \leftarrow *relations* + CLASSIFYRELATION(*e1*, *e2*)

Exemplos de Features

American Airlines, a unit of AMR, immediately matched the move, spokesman **Tim Wagner** said

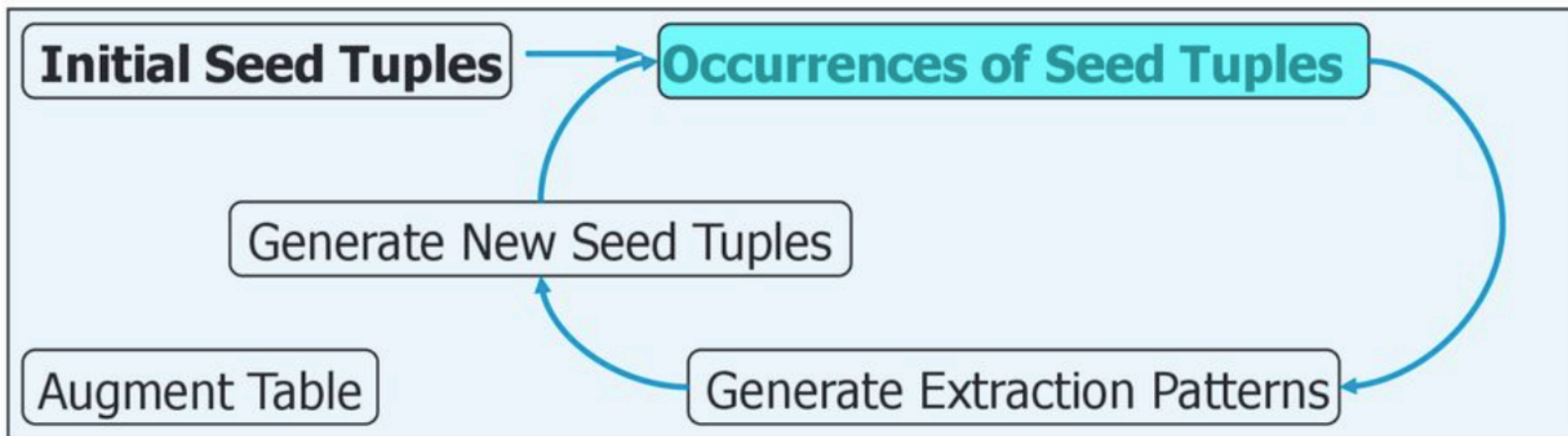
- Palavras ou bigramas ao redor
- Tipos das entidades
- Número de entidades entre as entidades candidatas
- Part-of-speech

Baseado em Rede Neural

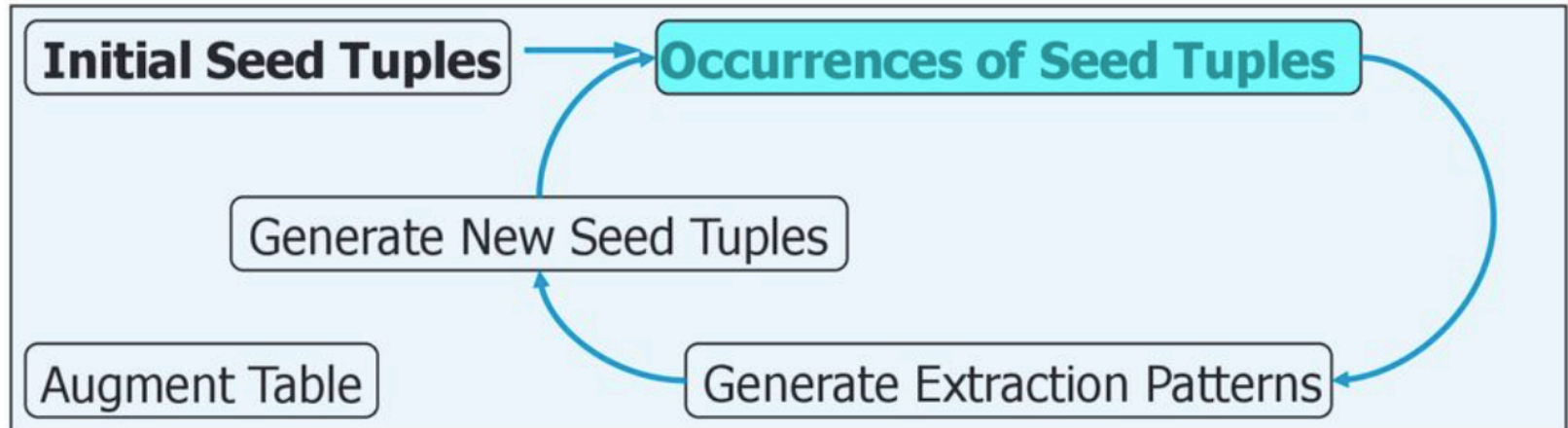


Bootstrapping

- Rótulos são difíceis de obter



Bootstrapping



Occurrences of
seed tuples:

ORGANIZATION	LOCATION
MICROSOFT	REDMOND
IBM	ARMONK
BOEING	SEATTLE
INTEL	SANTA CLARA

Computer servers at **Microsoft**'s headquarters in **Redmond**..

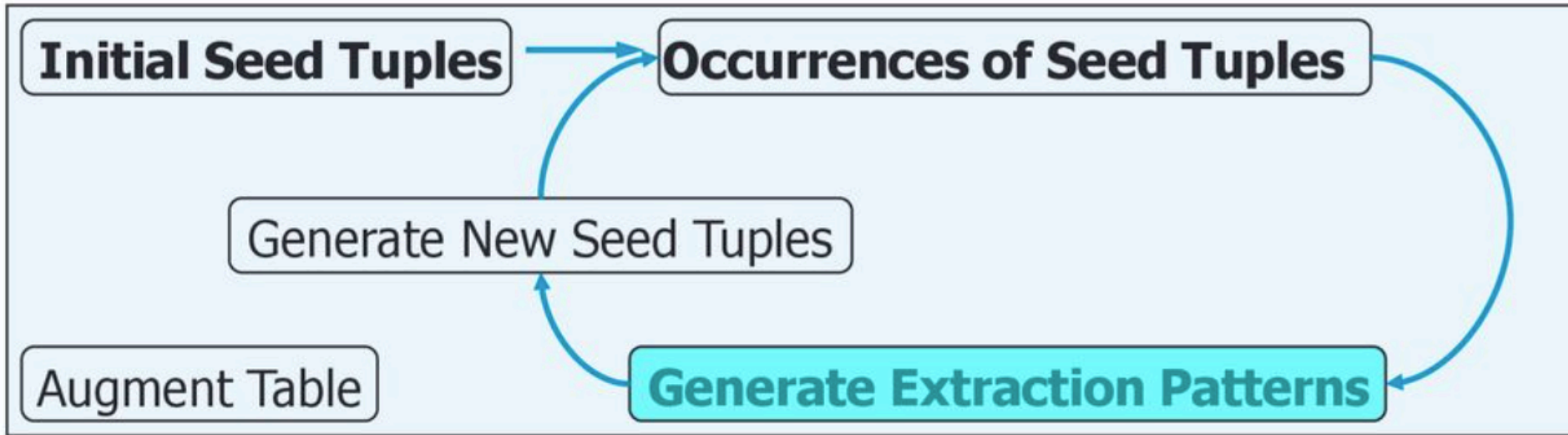
In mid-afternoon trading, share of **Redmond**-based **Microsoft** fell...

The **Armonk**-based **IBM** introduced a new line...

The combined company will operate from **Boeing**'s headquarters in **Seattle**.

Intel, **Santa Clara**, cut prices of its Pentium processor.

Bootstrapping



Bootstrapping: Gerando Novas Sementes

- Uso dos padrões aprendidos

ORGANIZATION	LOCATION
AG EDWARDS	ST LUIS
157TH STREET	MANHATTAN
7TH LEVEL	RICHARDSON
3COM CORP	SANTA CLARA
3DO	REDWOOD CITY
JELLIES	APPLE
MACWEEK	SAN FRANCISCO

Initial Seed Tuples

Occurrences of Seed Tuples

Generate New Seed Tuples

Augment Table

Generate Extraction Patterns

