

# Xử lý ngôn ngữ tự nhiên

## Chương 6: Parsing và Grammar

Khoa CNTT, Đại học Kỹ thuật - Công nghệ Cần Thơ  
Lưu hành nội bộ

# Nội dung

- 1 Parsing
- 2 Grammar

# Syntax

- Cú pháp (syntax) quy định sự sắp xếp các từ
  - Nhóm từ đi cùng nhau (constituency)
  - Quan hệ văn phạm (grammatical relations)
  - Các quy tắc lệ thuộc: thứ tự các từ loại

# Context-free grammar and trees

- Mô hình toán học mô tả cấu trúc ngôn ngữ
- Terminals và non-terminals
- Nguồn gốc (derivation)
- Parse tree
- Start symbol

*Noun* → *flight* | *breeze* | *trip* | *morning* | ...

*Verb* → *is* | *prefer* | *like* | *need* | *want* | *fly* ...

*Adjective* → *cheapest* | *non-stop* | *first* | *latest* | *other* | *direct* | ...

*Pronoun* → *me* | *I* | *you* | *it* | ...

*Proper-Noun* → *Alaska* | *Baltimore* | *Los Angeles* | *Chicago* | *United* | *American* | ...

*Determiner* → *the* | *a* | *an* | *this* | *these* | *that* | ...

*Preposition* → *from* | *to* | *on* | *near* | ...

*Conjunction* → *and* | *or* | *but* | ...

The lexicon for  $L_0$

$S \rightarrow NP VP$

$NP \rightarrow$  *Pronoun*

| *Proper-Noun*

| *Det Nominal*

$Nominal \rightarrow$  *Noun Nominal*

| *Noun*

$VP \rightarrow$  *Verb*

| *Verb NP*

| *Verb NP PP*

| *Verb PP*

$PP \rightarrow$  *Preposition NP*

I + want a morning flight

I

Los Angeles

a + flight

morning + flight

flights

do

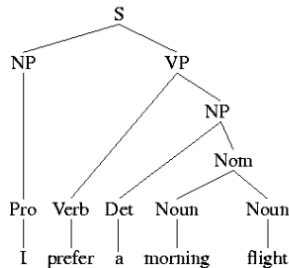
want + a flight

leave + Boston + in the morning

leaving + on Thursday

from + Los Angeles

The grammar for  $L_0$



- 1 Parsing
- 2 Grammar

# Parsing

- Parsing: sự phân tích ngữ pháp (từ, câu)
- Còn được gọi là syntactic analysis
- Là quá trình xác định xem một chuỗi các words liên tục nhau có đúng ngữ pháp hay không
- Chia nhỏ câu thành từ, ngữ và xác định nhóm thành phần: noun - danh từ, verb - động từ
- Tổng hợp hóa tác vụ POS

# Grammatically vs. semantically

- Ví dụ câu: John bought a book
  - Câu đúng ngữ pháp (grammatically correct)
- Ví dụ câu: book bought a John
  - Câu đúng ngữ pháp (grammatically correct)
  - Câu sai ngữ nghĩa (semantically incorrect)



# Nhóm Parsing

- Parsing được chia thành 2 nhóm:
  - Top-down parsing: bắt đầu từ ký tự bắt đầu (start symbol) và tiếp tục cho đến khi gặp các thành phần đơn (individual components)
  - Bottom-up parsing: bắt đầu từ các thành phần đơn và tiếp tục cho đến ký tự bắt đầu

# Treebank corpus

- Treebank là corpus gồm 199 documents được cung cấp kèm theo thư viện NLTK
- <https://en.wikipedia.org/wiki/Treebank>

```
1 import nltk
2 import nltk.corpus
3 from nltk.corpus import treebank
4
5 print(nltk.corpus.treebank.fileids())
6 print(len(nltk.corpus.treebank.fileids()))
7
8 #199
```

[illegible]

```
9
10 print(treebank.words('wsj_0007.mrg'))
11 print(len(treebank.words('wsj_0007.mrg')))
12 print(treebank.tagged_words('wsj_0007.mrg'))
13
14 #['McDermott', 'International', 'Inc.', 'said', '0', ...]
15 #75
16 #[( 'McDermott', 'NNP'), ('International', 'NNP'), ...]
```

# Penn Treebank corpus

```
1 import nltk
2 from nltk.corpus import treebank
3 print(treebank.sents('wsj_0007.mrg')[2])
4 print(treebank.parsed_sents('wsj_0007.mrg')[2])
5
6 #['Bailey', 'Controls', ',', 'based', '*', 'in', 'Wickliffe', ',', 'Ohio', ',', 'makes', 'computerized', 'industrial', 'controls', 'systems', '.']
```

```

(S
  (NP-SBJ
    (NP (NNP Bailey) (NNP Controls))
    (, ,)
    (VP
      (VBN based)
      (NP (-NONE- *))
      (PP-LOC-CLR
        (IN in)
        (NP (NP (NNP Wickliffe)) (, ,) (NP (NNP Ohio)))))
      (, ,))
  (VP
    (VBZ makes)
    (NP
      (JJ computerized)
      (JJ industrial)
      (NNS controls)
      (NNS systems)))
    (. .))

```

# POS tags in Penn Treebank

1	CC	Coordinating conjunction
2	CD	Cardinal number
3	DT	Determiner
4	EX	Existential there
5	FW	Foreign word
6	IN	Preposition or subordinating conjunction
7	JJ	Adjective
8	JJR	Adjective, comparative
9	JJS	Adjective, superlative
10	LS	List item marker

# POS tags in Penn Treebank

11	MD	Modal
12	NN	Noun, singular or mass
13	NNS	Noun, plural
14	NNP	Proper noun, singular
15	NNPS	Proper noun, plural
16	PDT	Predeterminer
17	POS	Possessive ending
18	PRP	Personal pronoun
19	PRP\$	Possessive pronoun
20	RB	Adverb



# POS tags in Penn Treebank

21	RBR	Adverb, comparative
22	RBS	Adverb, superlative
23	RP	Particle
24	SYM	Symbol
25	TO	<i>to</i>
26	UH	Interjection
27	VB	Verb, base form
28	VBD	Verb, past tense
29	VBG	Verb, gerund or present participle
30	VBN	Verb, past participle

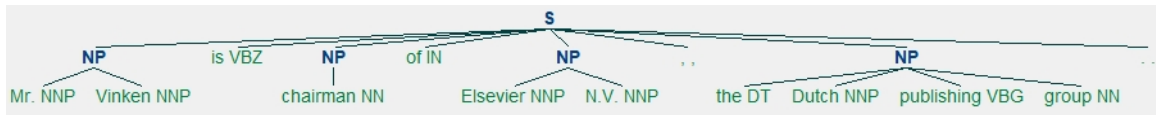
# POS tags in Penn Treebank

31	VBP	Verb, non-3rd person singular present
32	VBZ	Verb, 3rd person singular present
33	WDT	Wh-determiner
34	WP	Wh-pronoun
35	WP\$	Possessive wh-pronoun
36	WRB	Wh-adverb

```

1 import nltk
2 from nltk.corpus import treebank_chunk
3 print(treebank_chunk.chunked_sents()[1])
4 treebank_chunk.chunked_sents()[1].draw()

```



```

1 import nltk
2 from nltk.corpus import treebank_chunk
3
4 print(treebank_chunk.chunked_sents()[1])
5
6 print(treebank_chunk.chunked_sents()[1].leaves())
7 #(['Mr.', 'NNP'), ('Vinken', 'NNP'), ('is', 'VBZ'), ('chairman', 'NN'),
   ('of', 'IN'), ('Elsevier', 'NNP'), #('N.V.', 'NNP'), (',', ', ', ', '),
   ('the', 'DT'), ('Dutch', 'NNP'), ('publishing', 'VBG'), ('group', '
   NN'), (',', ', ')]
8
9 print(treebank_chunk.chunked_sents()[1].pos())
10 #(((('Mr.', 'NNP'), 'NP'), (((('Vinken', 'NNP'), 'NP'), (((('is', 'VBZ'), 'S
   '), (((('chairman', 'NN'), 'NP'), (((('of', 'IN'), 'S'), (((('Elsevier',
   'NNP'), 'NP'), (((('N.V.', 'NNP'), 'NP'), ((((', ', ', ', ', '), 'S'), (((('the
   ', 'DT'), 'NP'), #(((('Dutch', 'NNP'), 'NP'), (((('publishing', 'VBG'),
   'NP'), (((('group', 'NN'), 'NP'), ((((', ', ', '), 'S'))])

```

```

11
12 print(treebank_chunk.chunked_sents()[1].productions())
13
14 #[S -> NP ('is', 'VBZ') NP ('of', 'IN') NP (',', ',') NP (',', ','), NP
    -> ('Mr.', 'NNP') ('Vinken', 'NNP'), NP -> #('chairman', 'NN'), NP
    -> ('Elsevier', 'NNP') ('N.V.', 'NNP'), NP -> ('the', 'DT') ('
    Dutch', 'NNP') #('publishing', 'VBG') ('group', 'NN')]

```

1 Parsing

2 Grammar

# Context Free Grammar (CFG)

- Văn phạm phi ngữ cảnh (CFG) bao gồm các thành phần:
  - Một tập các nút không kết thúc (non terminal nodes) (N)
  - Một tập các nút kết thúc (terminal nodes) (T)
  - Ký hiệu bắt đầu (S)
  - Tập các luật dẫn xuất (production rules) (P) có dạng  $A \rightarrow a$

# Các luật CFG

- Luật cấu trúc ngữ (phrase structure rules) xác định thông qua các luật dẫn xuất  $A \rightarrow a$
- Luật cấu trúc câu (sentence structure rules)
  - Cấu trúc tường thuật (declarative structure): chủ ngữ - vị ngữ
  - Cấu trúc mệnh lệnh (imperative structure): câu ra lệnh, bắt đầu bằng động từ
  - Cấu trúc Yes-No
  - Cấu trúc câu hỏi Wh (Wh-question)



# Các luật CFG

- $S \rightarrow NP VP$
- $S \rightarrow VP$
- $S \rightarrow Aux NP VP$
- $S \rightarrow Wh-NP VP$
- $S \rightarrow Wh-NP Aux NP VP$
- $NP \rightarrow (Det) (AP) Nom (PP)$
- $VP \rightarrow Verb (NP) (NP) (PP)$

# Ví dụ CFG 1

```
1 import nltk
2 from nltk import Nonterminal, nonterminals, Production, CFG
3 nonterminal1 = Nonterminal('NP')
4 nonterminal2 = Nonterminal('VP')
5 nonterminal3 = Nonterminal('PP')
6
7 print(nonterminal1.symbol())
8 print(nonterminal2.symbol())
9 print(nonterminal3.symbol())
10 #NP
11 #VP
12 #PP
13
14 print(nonterminal1==nonterminal2)
15 print(nonterminal2==nonterminal3)
16 print(nonterminal1==nonterminal3)
```

```
17 #False
18 #False
19 #False
20
21 S, NP, VP, PP = nonterminals('S, NP, VP, PP')
22 N, V, P, DT = nonterminals('N, V, P, DT')
23 production1 = Production(S, [NP, VP])
24 production2 = Production(NP, [DT, NP])
25 production3 = Production(VP, [V, NP, NP, PP])
26
27 print(production1.lhs())
28 print(production1.rhs())
29 print(production3.lhs())
30 print(production3.rhs())
31 #S
32 #(NP, VP)
33 #VP
34 #(V, NP, NP, PP)
```

```
35  
36 print(production3 == Production(VP, [V,NP,NP,PP]))  
37 print(production2 == production3)  
38 #True  
39 #False
```

# Ví dụ CFG 2

```

1 from nltk import Nonterminal, nonterminals, Production, CFG
2 nt1 = Nonterminal('NP')
3 nt2 = Nonterminal('VP')
4 print(nt1.symbol())
5 print(nt1 == Nonterminal('NP'))
6 #NP
7 #True
8
9 S, NP, VP, PP = nonterminals('S, NP, VP, PP')
10 N, V, P, DT = nonterminals('N, V, P, DT')
11 prod1 = Production(S, [NP, VP])
12 prod2 = Production(NP, [DT, NP])
13 print(prod1.lhs())
14 print(prod1.rhs())
15 print(prod1 == Production(S, [NP, VP]))

```

```
16 #S
17 #(NP, VP)
18 #True
19
20 grammar = CFG.fromstring("""
21     S -> NP VP
22     PP -> P NP
23     NP -> 'the' N | N PP | 'the' N PP
24     VP -> V NP | V PP | V NP PP
25     N -> 'cat'
26     N -> 'dog'
27     N -> 'rug'
28     V -> 'chased'
29     V -> 'sat'
30     P -> 'in'
31     P -> 'on'
32 """)
33 print(grammar)
```

```
34 #Grammar with 15 productions (start state = S)
35 # S → NP VP
36 # PP → P NP
37 # NP → 'the' N
38 # NP → N PP
39 # NP → 'the' N PP
40 # VP → V NP
41 # VP → V PP
42 # VP → V NP PP
43 # N → 'cat'
44 # N → 'dog'
45 # N → 'rug'
46 # V → 'chased'
47 # V → 'sat'
48 # P → 'in'
49 # P → 'on'
```

# ATIS grammar

- Các luật dẫn xuất được cung cấp bởi các bộ văn phạm

```
1 import nltk
2 gram1 = nltk.data.load('grammars/large_grammars/atis.cfg')
3 print(grammar1)
4
5 #Grammar with 5517 productions (start state = SIGMA)
6 #ABBCL_NP -> QUANP_DTI QUANP_DTI QUANP_CD AJP_JJ NOUN_NP PRPRTCL_VBG
7 #ADJ_ABL -> only
8 #ADJ_ABL -> such
9 #ADJ_AP -> pt_adj_ap
10 #ADJ_AP -> other
11 #ADJ_AP -> last
12 #ADJ_AP -> less
13 #ADJ_AT -> the
14 #...
```



# Trích xuất câu trong ATIS grammar

```
1 import nltk
2 sent = nltk.data.load('grammars/large_grammars/atis_sentences.txt')
3 sent = nltk.parse.util.extract_test_sentences(sent)
4 print(len(sent))
5 testingsent=sent[25]
6 print(testingsent)
7 print(testingsent[0])
8 print(testingsent[1])
9
10 #98
11 #(['list', 'those', 'flights', 'that', 'stop', 'over', 'in', 'salt', '
    lake', 'city', '.'], 11)
12 #(['list', 'those', 'flights', 'that', 'stop', 'over', 'in', 'salt', '
    lake', 'city', '.']
13 #11
```

# Recursive Descent Parser

```
1 from nltk import CFG
2 from nltk.parse import RecursiveDescentParser
3
4 grammar = CFG.fromstring("""
5     S -> NP VP
6     PP -> P NP
7     NP -> 'the' N | N PP | 'the' N PP
8     VP -> V NP | V PP | V NP PP
9     N -> 'cat'
10    N -> 'dog'
11    N -> 'rug'
12    V -> 'chased'
13    V -> 'sat'
14    P -> 'in'
15    P -> 'on'
16 """)
```

```

17
18 rd = RecursiveDescentParser(grammar)
19
20 sentence1 = 'the cat chased the dog'.split()
21 sentence2 = 'the cat chased the dog on the rug'.split()
22
23 for t in rd.parse(sentence1):
24     print(t)
25 for t in rd.parse(sentence2):
26     print(t)
27
28 #(S (NP the (N cat)) (VP (V chased) (NP the (N dog))))
29 #(S
30 #(NP the (N cat))
31 #(VP (V chased) (NP the (N dog) (PP (P on) (NP the (N rug))))))
32 #(S
33 #(NP the (N cat))
34 #(VP (V chased) (NP the (N dog)) (PP (P on) (NP the (N rug))))

```

# Chart Parser

```
1 import nltk
2
3 sent = 'I saw a dog'
4 nltk.parse.chart.demo(2, print_times=False, trace=1, sent=sent,
5                        numparses=1)
6
7 #* Sentence:
8 #I saw a dog
9 #[ 'I', 'saw', 'a', 'dog' ]
10 #
11 #* Strategy: Bottom-up
```

.	I	.	saw	.	a	.	dog	.	
	[-----]	.		.		.		.	[0:1] 'I'
.		[-----]		.		.		.	[1:2] 'saw'
.		.	[-----]			.		.	[2:3] 'a'
.		.	.	[-----]				.	[3:4] 'dog'
>	.	.	.	.		.		.	[0:0] NP -> * 'I'
	[-----]	.	.	.		.		.	[0:1] NP -> 'I' *
>	.	.	.	.		.		.	[0:0] S -> * NP VP
>	.	.	.	.		.		.	[0:0] NP -> * NP PP
	[----->	.	.	.		.		.	[0:1] S -> NP * VP
	[----->	.	.	.		.		.	[0:1] NP -> NP * PP
.	>	.	.	.		.		.	[1:1] Verb -> * 'saw'
.	[-----]	.	.	.		.		.	[1:2] Verb -> 'saw' *
.	>	.	.	.		.		.	[1:1] VP -> * Verb NP
.	>	.	.	.		.		.	[1:1] VP -> * Verb
.	[----->	.	.	.		.		.	[1:2] VP -> Verb * NP
.	[-----]	.	.	.		.		.	[1:2] VP -> Verb *
.	>	.	.	.		.		.	[1:1] VP -> * VP PP
	[-----]	.	.	.		.		.	[0:2] S -> NP VP *
.	[----->	.	.	.		.		.	[1:2] VP -> VP * PP

```
1 import nltk
2 sent = 'I saw John with a dog'
3 nltk.parse.chart.demo(1, print_times=False, trace=0, sent=sent,
   numparses=2)
4
5 #* Sentence:
6 #I saw John with a dog
7 #['I', 'saw', 'John', 'with', 'a', 'dog']
8 #
9 #* Strategy: Top-down
10 #
11 #Nr edges in chart: 48
12 #(S
13 #(NP I)
14 #(VP (Verb saw) (NP (NP John) (PP with (NP (Det a) (Noun dog))))))
15 #(S
16 #(NP I)
17 #(VP (VP (Verb saw) (NP John)) (PP with (NP (Det a) (Noun dog))))))
```

```
1 import nltk
2 sent = 'I saw John with a dog'
3 nltk.parse.chart.demo(5, print_times=False, trace=1, sent=sent,
4                        numparses=2)
5
6 #* Sentence:
7 #I saw John with a dog
8 #['I', 'saw', 'John', 'with', 'a', 'dog']
9 #
10 #* Strategy: Stepping (top-down vs bottom-up)
```

\*\*\* SWITCH TO TOP DOWN

```

| [-----] . . . . | [0:1] 'I'
| . [-----] . . . | [1:2] 'saw'
| . . [-----] . . | [2:3] 'John'
| . . . [-----] . | [3:4] 'with'
| . . . . [-----] | [4:5] 'a'
| . . . . . [-----] | [5:6] 'dog'
|> . . . . . | [0:0] S -> * NP VP
|> . . . . . | [0:0] NP -> * NP PP
|> . . . . . | [0:0] NP -> * Det Noun
|> . . . . . | [0:0] NP -> * 'I'
| [-----] . . . . | [0:1] NP -> 'I' *
| [-----> . . . . | [0:1] S -> NP * VP
| [-----> . . . . | [0:1] NP -> NP * PP
| . > . . . . | [1:1] VP -> * VP PP
| . > . . . . | [1:1] VP -> * Verb NP
| . > . . . . | [1:1] VP -> * Verb
| . > . . . . | [1:1] Verb -> * 'saw'
| . [-----] . . . | [1:2] Verb -> 'saw' *
| . [-----> . . . | [1:2] VP -> Verb * NP
| . [-----> . . . | [1:2] VP -> Verb * NP

```



# Probabilistic CFG

- Là dạng cú pháp CFG có thêm thông tin về xác suất xảy ra luật

```
1 import nltk
2 from nltk.corpus import treebank
3 from itertools import islice
4 from nltk.grammar import PCFG, induce_pcfg, toy_pcfg1, toy_pcfg2
5
6 gram2 = PCFG.fromstring("""
7     A -> B B [.3] | C B C [.7]
8     B -> B D [.5] | C [.5]
9     C -> 'a' [.1] | 'b' [0.9]
10    D -> 'b' [1.0]
11 """)
```

```

12 prod1 = gram2 productions() [0]
13 print(prod1)
14 prod2 = gram2 productions() [1]
15 print(prod2)
16 #A -> B B [0.3]
17 #A -> C B C [0.7]
18
19 print(prod2.lhs())
20 print(prod2.rhs())
21 print((prod2.prob()))
22 print(gram2.start())
23 print(gram2 productions())
24
25 #A
26 #(C, B, C)
27 #0.7
28 #A
29 #[A -> B B [0.3], A -> C B C [0.7], B -> B D [0.5], B -> C [0.5], C ->
    'a' [0.1], C -> 'b' [0.9], D -> 'b' [1.0]]

```

# Kết hợp grammar và parser

```
1 import nltk
2 from nltk.corpus import treebank
3 from itertools import islice
4 from nltk.grammar import PCFG, induce_pcfg, toy_pcfg1, toy_pcfg2
5 grammar = toy_pcfg2
6 print(grammar)
7
8 #Grammar with 23 productions (start state = S)
9 #S -> NP VP [1.0]
10 #VP -> V NP [0.59]
11 #VP -> V [0.4]
12 #VP -> VP PP [0.01]
13 #NP -> Det N [0.41]
14 #NP -> Name [0.28]
15 #NP -> NP PP [0.31]
16 #PP -> P NP [1.0]
```

```
17 from nltk.parse import pchart
18
19 tokens = "Jack saw Bob with my cookie".split()
20 parser = pchart.InsideChartParser(grammar)
21 for t in parser.parse(tokens):
22     print(t)
23 #(S
24 #(NP (Name Jack))
25 #(VP
26 #(V saw)
27 #(NP
28 #(NP (Name Bob))
29 #(PP (P with) (NP (Det my) (N cookie)))))) (p=6.31607e-06)
30 #(S
31 #(NP (Name Jack))
32 #(VP
33 #(VP (V saw) (NP (Name Bob)))
34 #(PP (P with) (NP (Det my) (N cookie)))))) (p=2.03744e-07)
```

# Visualization

- <https://demos.explosion.ai/displacy/>

