



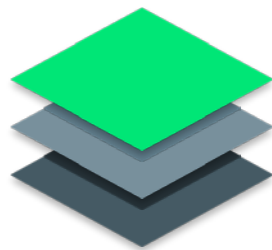
Dependency Parsing of Korean using SyntaxNet

고려대학교 컴퓨터학과

Andrew Matteson

Web : www.andrewmatteson.name

E-mail : amatteson@korea.ac.kr



Deep Learning for NLP lecture

고려대학교 정보대학 컴퓨터학과

2017

Index

3

1

Introduction

Dependency
Grammar

SyntaxNet

2

Challenges

SyntaxNet

Korean

3

Preparation

Corpus

SyntaxNet

POS Tagger

Custom Scripts

4

Training

Configuration

Training Script

Greedy Training

Structured Training

5

Execution

Workflow

Tree Visualization

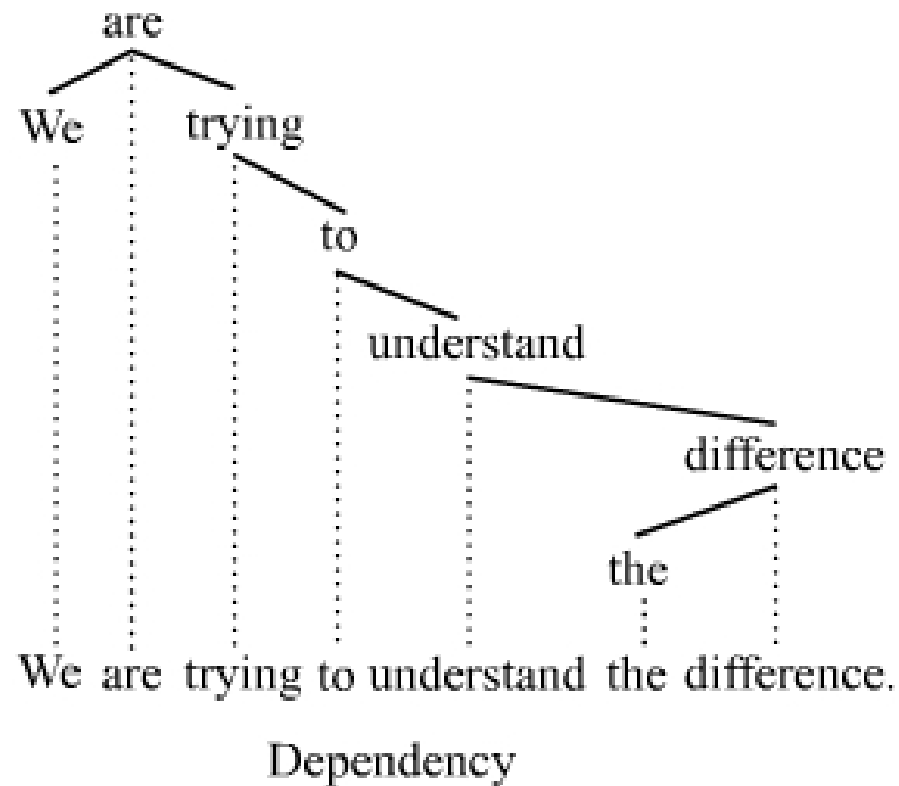
Example

Introduction

The background features a series of overlapping chevron shapes pointing to the right. The chevrons are in various shades of blue, ranging from a very light, almost white blue on the far right to a dark navy blue on the left. A prominent, bright green line runs diagonally across the image, starting from the top center and extending towards the bottom right, passing over the chevrons.

Dependency Grammar

- Analyze dependencies between certain words



SyntaxNet (1/3)

- Package by Google allowing:
 - Tokenization
 - Morphological analysis
 - POS tagging
 - Dependency parsing
- Only supports Python 2

SyntaxNet (2/3)

7

- Embeds sets of features deeply
- Hidden layer of words can also serve as word vector representation
- Two network graph structures are provided:
 - **Greedy:** picks only best #1 choice at each step of parsing
 - **Structured:** allows beam search with Conditional Random Field objective over several hypotheses

SyntaxNet (3/3)

8

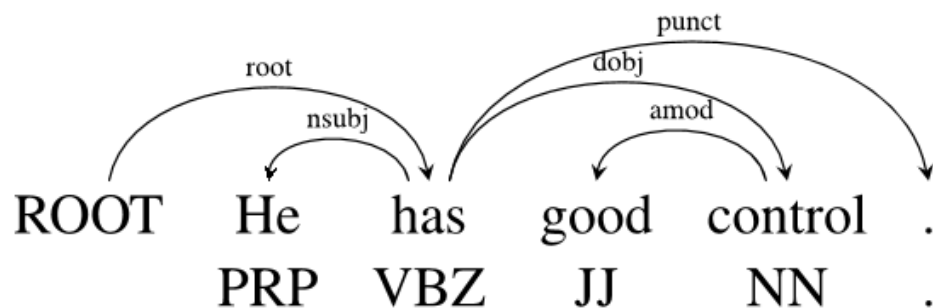
- Information typically flows from top (tokenizer) to bottom (dependency parser)
 - Output of tagging used as input to parsing
- “Parsey McParseface”

SyntaxNet Parser

- Arc-standard transition system
- Input sentences
 - Represented as a sequence of transition actions and labels
- Train input parser state against gold standard output parser state

Arc-standard Transition System

10



Correct transition: SHIFT



Transition	Stack	Buffer	A
	[ROOT]	[He has good control .]	\emptyset
SHIFT	[ROOT He]	[has good control .]	
SHIFT	[ROOT He has]	[good control .]	
LEFT-ARC (nsubj)	[ROOT has]	[good control .]	$A \cup \text{nsubj}(\text{has}, \text{He})$
SHIFT	[ROOT has good]	[control .]	
SHIFT	[ROOT has good control]	[.]	
LEFT-ARC (amod)	[ROOT has control]	[.]	$A \cup \text{amod}(\text{control}, \text{good})$
RIGHT-ARC (dobj)	[ROOT has]	[.]	$A \cup \text{dobj}(\text{has}, \text{control})$
...
RIGHT-ARC (root)	[ROOT]	[]	$A \cup \text{root}(\text{ROOT}, \text{has})$

Challenges

The background features a series of overlapping chevron shapes pointing to the right. The chevrons are in various shades of blue, ranging from a dark charcoal blue on the left to a very light, almost white blue on the right. A prominent, thick, bright green line runs diagonally across the image, starting from the top center and extending towards the bottom right, passing over the chevrons.

Challenges of SyntaxNet

12

- Lack of documentation
- Input features to parser are important
- Processes sentences as a series of tokens
 - Tokens assumed to be separated by whitespace
 - Separating Korean into tokens is difficult
- Can't yet use SyntaxNet tagger as input to parser in Korean case
 - Tagging Korean in SyntaxNet difficult due to lack of Korean features
- No GPU support (slow)

Challenges of Korean

13

- Morphologically complex
- Some markers in Korean are optional
 - 나는 닭발을 잘 안 먹어
 - 나 닭발 잘 안 먹어
- Honorifics/verb stem changes
 - 저는 닭발을 잘 안 먹어요
 - 우리 교수님은 닭발을 잘 안 드세요



Preparation

Prerequisites

15

- Sejong Written Language Corpus (Constituency-Tree Format)
 - “현대문어 구문분석 말뭉치” available on sejong.or.kr
 - Preprocessing/output to CoNLL-U Format
- Google SyntaxNet
 - Installation process is well-documented
- Third-party POS tagger
 - Not using SyntaxNet for tagging
- Custom scripts
 - Convert between necessary formats, etc...

Sejong Corpus (1/5)

16

- Available freely online under Miscellaneous section of sejong.or.kr
- Sample sentence:

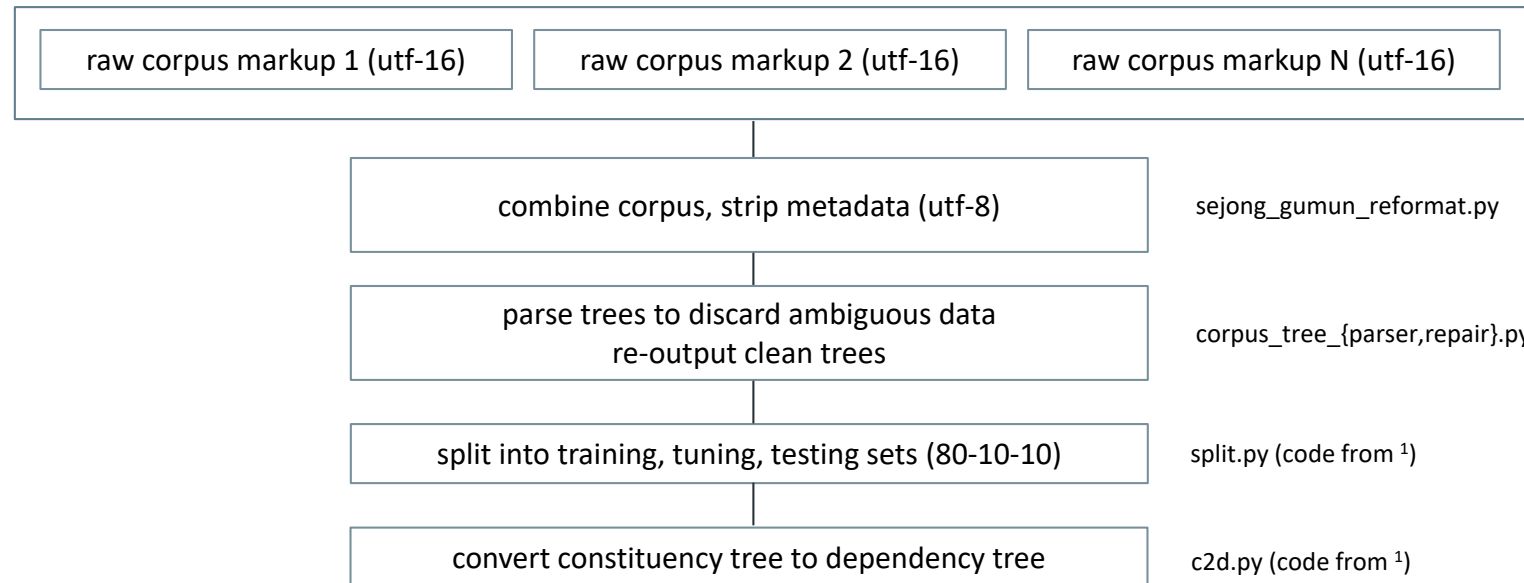
; 자신도 고교를 졸업하고 용돈을 벌 용량으로 띄엄띄엄 공사판을 돌아다니다 배관 기술을 익혀 눌러 앉았다는 전씨의 월수입은 1백 40여 만원.

(S (NP_SBJ (NP_MOD (S_MOD (NP_SBJ 자신/NNG + 도/JX) (VP_MOD (VP (VP (VP (NP_OBJ 고교/NNG + 를/JKO) (VP 졸업/NNG + 하/XSV + 고/EC)) (VP (NP_AJT (VP_MOD (NP_OBJ 용돈/NNG + 을/JKO) (VP_MOD 벌/VV + 께/ETM)) (NP_AJT 용량/NNG + 으로/JKB)) (VP (AP 띄엄띄엄/MAG) (VP (NP_OBJ 공사판/NNG + 을/JKO) (VP 돌아다니/VV + 다/EC)))))) (VP (NP_OBJ (NP 배관/NNG) (NP_OBJ 기술/NNG + 을/JKO)) (VP 익히/VV + 어/EC))) (VP_MOD (VP 누르/VV + 어/EC) (VP_MOD 앉/VV + 앉/EP + 다는/ETM)))) (NP_MOD 전/NNP + 씨/NNB + 의/JKG)) (NP_SBJ 월/NNG + 수입/NNG + 은/JX)) (NP (NP 1/SN + 백/NR) (NP 40/SN + 여/XSN)) (NP 만/NR + 원/NNG + ./SF)))

Sejong Corpus (2/5)

17

Overview



¹ <https://github.com/dsindex/syntaxnet>

Sejong Corpus (3/5)

18

Result

```
<!DOCTYPE tei.2 SYSTEM "c:\sgml\dtd\tei2.dtd" [  
  <!ENTITY % TEI.corpus "INCLUDE">  
  <!ENTITY % TEI.extensions.ent SYSTEM "sejong1.ent">  
  <!ENTITY % TEI.extensions.dtd SYSTEM "sejong1.dtd">  
  
<tei.2>  
<teiHeader>  
  <fileDesc>  
    <titleStmt>  
      <title>조선일보 생활(93), 구문 분석 전자파일</title>  
      <author>조선일보사</author>  
      <sponsor>대한민국 문화관광부</sponsor>  
      <respStmt>  
        <resp>문헌입력, 표준화, 구문 정보 부착</resp>  
      ...  
    <body>  
      ; 1993/06/08 19  
(NP      (NP 1993/SN + //SP + 06/SN + //SP + 08/SN)  
          (NP 19/SN))
```



```
; 1993/06/08 19  
(NP      (NP 1993/SN + //SP + 06/SN + //SP + 08/SN)  
          (NP 19/SN))  
  
; 엠마누엘 웅가로 /  
(NP      (NP      (NP 엠마누엘/NNP)  
                  (NP 웅가로/NNP))  
          (X //SP))  
  
; 의상서 실내 장식품으로...  
(NP_AJT      (NP_AJT 의상/NNG + 서/JKB)  
              (NP_AJT      (NP 실내/NNG)  
                            (NP_AJT 장식품/NNG + 으로/JKB + .../SE)))  
  
; 디자인 세계 넓혀  
(VP      (NP_OBJ      (NP 디자인/NNG)  
                      (NP_OBJ 세계/NNG))  
          (VP 넓히/VV + 어/EC))
```

Sejong Corpus (4/5)

19

- Depth-first traversal of constituency tree
- Output to CoNLL-U format (introduced later)
- Determine HEAD of all leaf nodes

Sejong Corpus (5/5)

20

Constituency Tree -> Dependency Tree (c2d)

; 프랑스의 세계적인 의상 디자이너 엠마누엘 웅가로가 실내 장식용 식물 디자인
 11로 나눴다.

```
(S
  (NP_SBJ
    (NP
      (NP_MOD 프랑스/NNP + 의/JKG)
      (NP
        (VNP_MOD 세계/NNG + 적/XSN + 이/VCP + ㄴ/ETM)
        (NP
          (NP 의상/NNG)
          (NP 디자이너/NNG)))
        (NP_SBJ
          (NP 엠마누엘/NNP)
          (NP_SBJ 웅가로/NNP + 가/JKS)))
        (VP
          (NP_AJT
            (NP
              (NP 실내/NNG)
              (NP 장식/NNG + 용/XSN))
              (NP 식물/NNG))
            (NP_AJT 디자이너/NNG + 로/JKB))
            (VP 나서/VV + 었/EP + 다/EF + ./SF)))
  )
)
```



1	프랑스의	프랑스/NNP + 의/JKG	NP_MOD	4	
2	세계적인	세계/NNG + 적/XSN + 이/VCP + ㄴ/ETM			VNP_MOD 4
3	의상	의상/NNG	NP	4	
4	디자이너	디자이너/NNG	NP	6	
5	엠마누엘	엠마누엘/NNP	NP	6	
6	웅가로가	웅가로/NNP + 가/JKS	NP_SBJ	11	
7	실내	실내/NNG	NP	8	
8	장식용	장식/NNG + 용/XSN	NP	9	
9	식물	식물/NNG	NP	10	
10	디자이너로	디자이너/NNG + 로/JKB	NP_AJT	11	
11	나눴다. 나서/VV + 었/EP + 다/EF + ./SF	VP		0	

Credits to Myungchul Shin¹

¹ <https://github.com/dsindex/syntaxnet>

Google SyntaxNet

21

- <https://github.com/tensorflow/models/tree/master/syntaxnet>
- As of writing, Python 2.7 only
- Install prerequisites for SyntaxNet
- Compilation process hogs huge amount of RAM and CPU
- Compile and verify SyntaxNet tests succeed

POS Tagger (1/2)

22

- Reassembly requires POS tagger that supports 어절 indexing
 - Reassembly makes output human-readable
 - Provides index of input character in morphological analysis
- Unaware of open source tagger providing this
 - Except Komoran 3.0 beta

POS Tagger (2/2)

23

- Komoran 3.0 Beta
 - <http://shineware.tistory.com/entry/KOMORAN-30-beta>
 - <https://github.com/shin285/KOMORAN>
- Requires dependency jar files by same author
- Compile in Eclipse Java workspace

Custom Scripts

24

- Myungchul Shin's github
 - <https://github.com/dsindex/syntaxnet>
 - split.py : 80/10/10 training/tuning/testing
 - c2d.py : constituency tree -> dependency tree
- Remove erroneous sentences from corpus
 - Until c2d.py succeeds
 - Manually or write script
- Reassembly of components to Eojeol (optional)
 - Use token 어절 index from POS tagger
- Tree visualization
 - CoNLL-U -> tree
 - SVG file output

Additional Info

25

- Read available documentation
- Google SyntaxNet
 - Training process/etc.
 - <https://github.com/tensorflow/models/tree/master/syntaxnet>
- Myungchul Shin's github
 - Sejong corpus related info
 - <https://github.com/dsindex/syntaxnet>



Training

SyntaxNet Configuration (1/2)

27

- Input corpus needs to be in CoNLL-U Format

ID	FORM	LEMMA	UPOSTAG	XPOSTAG	FEATS	HEAD	DEPREL	DEPS	MISC
1	나	_	NP	NP	_	2	MOD	_	_
2	는	_	JX	JX	_	5	NP_SBJ	_	_
3	자연언 어처리	_	NNP	NNP	_	4	MOD	_	_
4	를	_	JKO	JKO	_	5	NP_OBJ	_	_
5	좋아하	_	VV	VV	_	6	MOD	_	_
6	ㄴ 다	_	EF	EF	_	7	MOD	_	_
7	.	_	SF	SF	_	0	ROOT	_	_

SyntaxNet Configuration (2/2)

28

- Configuration file called context.pbtxt
- Configures parser
- Designates input features to parser

This Model's Features (1/3)

29

`token.word`: word embedding for specified token

Feature Group: words (dim=64)	
<code>input.word</code>	<code>stack(1).child(1).word</code>
<code>input(1).word</code>	<code>stack(1).child(1).sibling(-1).word</code>
<code>input(2).word</code>	<code>stack(1).child(-1).word</code>
<code>input(3).word</code>	<code>stack(1).child(-1).sibling(1).word</code>
<code>stack.word</code>	<code>stack.child(2).word</code>
<code>stack(1).word</code>	<code>stack.child(-2).word</code>
<code>stack(2).word</code>	<code>stack(1).child(2).word</code>
<code>stack(3).word</code>	<code>stack(1).child(-2).word</code>
<code>stack.child(1).word</code>	
<code>stack.child(1).sibling(-1).word</code>	
<code>stack.child(-1).word</code>	
<code>stack.child(-1).sibling(1).word</code>	

This Model's Features (2/3)

`token.tag`: POS tag used in specified token (previous, current, or future tokens)

Feature Group: tags (dim=32)	
input.tag	stack(1).child(1).tag
input(1).tag	stack(1).child(1).sibling(-1).tag
input(2).tag	stack(1).child(-1).tag
input(3).tag	stack(1).child(-1).sibling(1).tag
stack.tag	stack.child(2).tag
stack(1).tag	stack.child(-2).tag
stack(2).tag	stack(1).child(2).tag
stack(3).tag	stack(1).child(-2).tag
stack.child(1).tag	
stack.child(1).sibling(-1).tag	
stack.child(-1).tag	
stack.child(-1).sibling(1).tag	

This Model's Features (3/3)

31

`token.label`: predicted DEPREL label (stack (previously processed) tokens only)

Feature Group: labels (dim=32)

`stack.child(1).label`

`stack.child(1).sibling(-1).label`

`stack.child(-1).label`

`stack.child(-1).sibling(1).label`

`stack(1).child(1).label`

`stack(1).child(1).sibling(-1).label`

`stack(1).child(-1).label`

`stack(1).child(-1).sibling(1).label`

`stack.child(2).label`

`stack.child(-2).label`

`stack(1).child(2).label`

`stack(1).child(-2).label`

SyntaxNet Training Script

32

Training Script

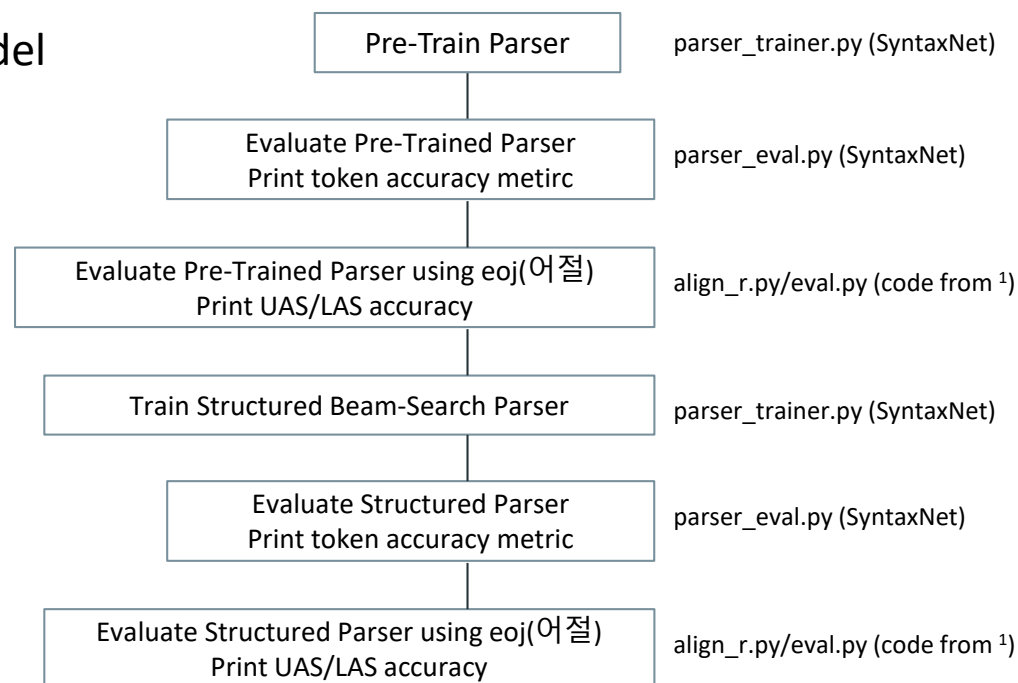
Bash script to train SyntaxNet model

```
GP_PARAMS=${HIDDEN_LAYER_PARAMS}-0.02-100-0.9
function train_parser {
    ${BINDIR}/parser_trainer \
    --arg_prefix=brain_parser \
    --batch_size=${BATCH_SIZE} \
    --compute_lexicon \
    --decay_steps=100 \
    --graph_builder=structured \
    --hidden_layer_sizes=${HIDDEN_LAYER_SIZES} \
    --learning_rate=0.02 \
    --momentum=0.9 \
    --beam_size=${BEAM_SIZE} \
    --output_path=${TMP_DIR} \
    --task_context=${TMP_DIR}/brain_parser/greedy/${LP} \
    --training_corpus=projectivized-training-corpus \
    --tuning_corpus=tagged-tuning-corpus \
    --params=${GP_PARAMS} \
    --pretrained_params=${TMP_DIR}/brain_parser/greedy/s \
    --pretrained_params_names=embedding_matrix_0,embedd \
    --num_epochs=10 \
    --report_every=25 \
    --checkpoint_every=200 \
    --logtostderr
}
```

```
...
pretrain_parser
evaluate_pretrained_parser
evaluate_pretrained_parser_by_eoj
train_parser
evaluate_parser
evaluate_parser_by_eoj
copy_model
```

Label Attachment Score (LAS) – % of tokens with correct HEAD and DEPREL

Unlabeled Attachment Score (UAS) – % of tokens with correct HEAD



¹ <https://github.com/dsindex/syntaxnet>

Greedy Training

33

- Required first step
- Usable model can be generated in just a couple hours

Greedy Training - Hyperparameters

34

- Batch size: 256
- Hidden layers:
 - 512 (ReLU activation)
 - 512 (ReLU activation)
- Output layer:
 - Softmax activation for choosing best transition action
- Objective: Minimize binary cross-entropy
- Optimizer: `tf.train.MomentumOptimizer`
 - Learning rate: 0.08
 - Momentum: 0.85
- Decay steps: 4400
- Number of epochs: 20

[number of individual tokens (incl. all features) input at once]

[decay learning rate by 0.96 every n steps]

[number of full passes through entire training set]

Structured Training

35

- Requires greedy model
- Uses beam search
- Identifies and fixes training faults in greedy model
- Can take a day or two

Structured Training - Hyperparameters

36

- Batch size: 256
- Hidden layers:
 - 512 (ReLU activation)
 - 512 (ReLU activation)
- Output layer:
 - Softmax activation for choosing best transition action
 - Beam search also performed
- Objective: Minimize binary cross-entropy
- Optimizer: `tf.train.MomentumOptimizer`
 - Learning rate: 0.02
 - Momentum: 0.90
- Beam size: 16
- Decay steps: 100
- Number of epochs: 10

Accuracy Assessment

37

- **Label Attachment Score (LAS)** - % of tokens with correct HEAD and DEPREL
- **Unlabeled Attachment Score (UAS)** - % of tokens with correct HEAD

Evaluated Sentence Count	
Training	Testing
19,569	3,788

	UAS		LAS	
	Training	Testing	Training	Testing
Greedy	92.17%	87.71%	89.56%	84.58%
Structured	94.61%	88.68%	92.13%	85.13%



Execution

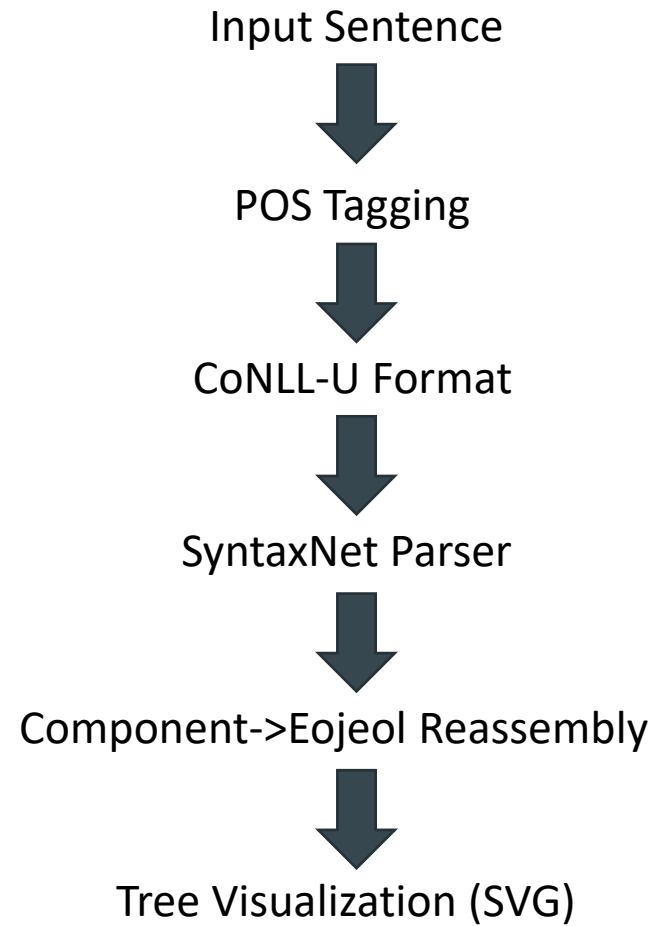
Execution

39

- So far, we have a trained model
- But running model is a whole different process

Workflow

40



POS Tagging (1/2)

41

- Invoke Komoran 3 on input sentence (Java)
- Save 어절 index of each token for later reassembly
 - Can be saved in CoNLL-U metadata field or externally
- Output to CoNLL-U Format
 - Only fill in POS tagging fields
 - HEAD or DEPREL later filled by SyntaxNet parser

POS Tagging (2/2)

42

- Input: “... 회사로 분할하기로 한 것은 ...”
- Part of speech tagging

ID	FORM	LEMMA	UPOSTAG	XPOSTAG	FEATS	HEAD	DEPREL	DEPS
11	회사	회사	NNG	NNG	—	—	—	—
12	로	로	JKB	JKB	—	—	—	—
13	분할	분할	NNG	NNG	—	—	—	—
14	하	하	XSV	XSV	—	—	—	—
15	기	기	ETN	ETN	—	—	—	—
16	로	로	JKB	JKB	—	—	—	—
17	한	한	MM	MM	—	—	—	—
18	것	것	NNB	NNB	—	—	—	—
19	은	은	JX	JX	—	—	—	—

회사	NNG	로	JKB	분할	NNG	하	XSV	기	ETN
로	JKB	한	MM	것	NNB	은	JX		



CoNLL-U Format



SyntaxNet Parser

SyntaxNet Parser (1/2)

44

- Parser CoNLL-U Output

HEAD (parent node ID)

ID	FORM	LEMMA	UPOSTAG	XPOSTAG	FEATS	HEAD (parent node ID)	DEPREL	DEPS
11	회사	—	NNG	NNG	—	12	MOD	—
12	로	—	JKB	JKB	—	13	NP_AJT	—
13	분할	—	NNG	NNG	—	14	MOD	—
14	하	—	XSV	XSV	—	15	MOD	—
15	기	—	ETN	ETN	—	16	MOD	—
16	로	—	JKB	JKB	—	60	VP_AJT	—
17	한	—	MM	MM	—	18	DP	—
18	것	—	NNB	NNB	—	19	MOD	—
19	은	—	JX	JX	—	60	NP_SBJ	—

Modifier

Eojeol Group

Eojeol Reassembly

45

- Component-→Eojeol Reassembly
 - Recombine prior MOD components into final Eojeol
 - Remap all tree node IDs

ID	FORM	LEMMA	HEAD	DEPREL
6	회사로	회사/NNG + 로/JKB	7	NP_AJT
7	분할하기로	분할/NNG + 하/XSV + 기 /ETN + 로/JKB	27	VP_AJT
8	한	한/MM	9	DP
9	것은	것/NNB + 은/JX	27	NP_SBJ

SyntaxNet Parser



Component->Eojeol Reassembly



Tree Visualization (SVG)

Tree Visualization (1/3)

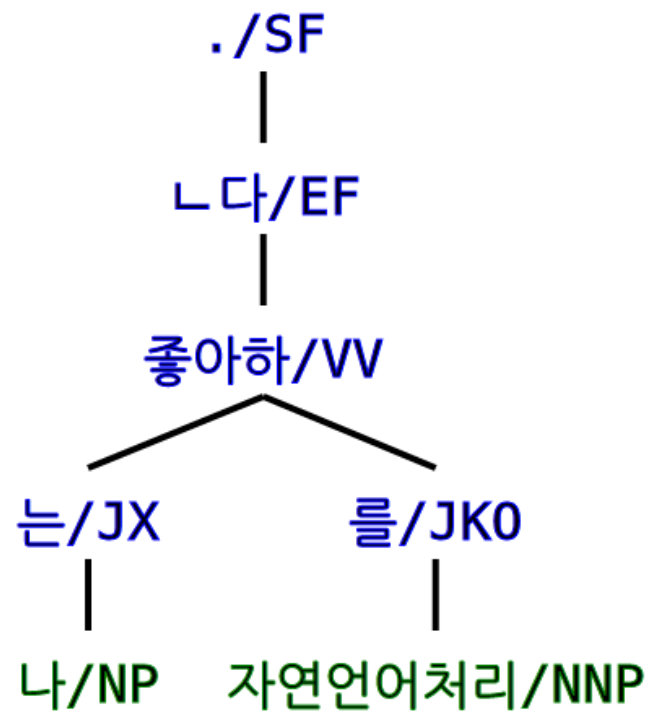
47

Given CoNLL-U format, we can easily generate a tree.

```
childList = {}
for elem in conllTable:
    if not(elem.HEAD in childList):
        childList[elem.HEAD] = []
    childList[elem.HEAD].append(elem)
    if(elem.HEAD == 0):
        rootElem = elem

def makeNode(elem, childList):
    node = Node()
    if elem.ID in childList:
        for child in childList[elem.ID]:
            node.children.append(makeNode(child, childList))
    return node

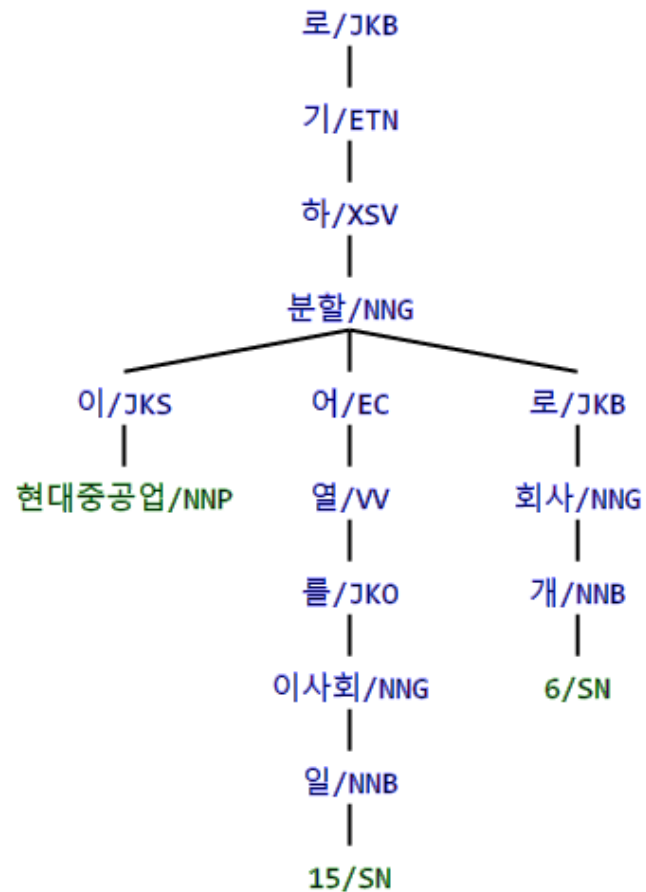
rootNode = makeNode(rootElem, childList)
```



Tree Visualization (2/3)

48

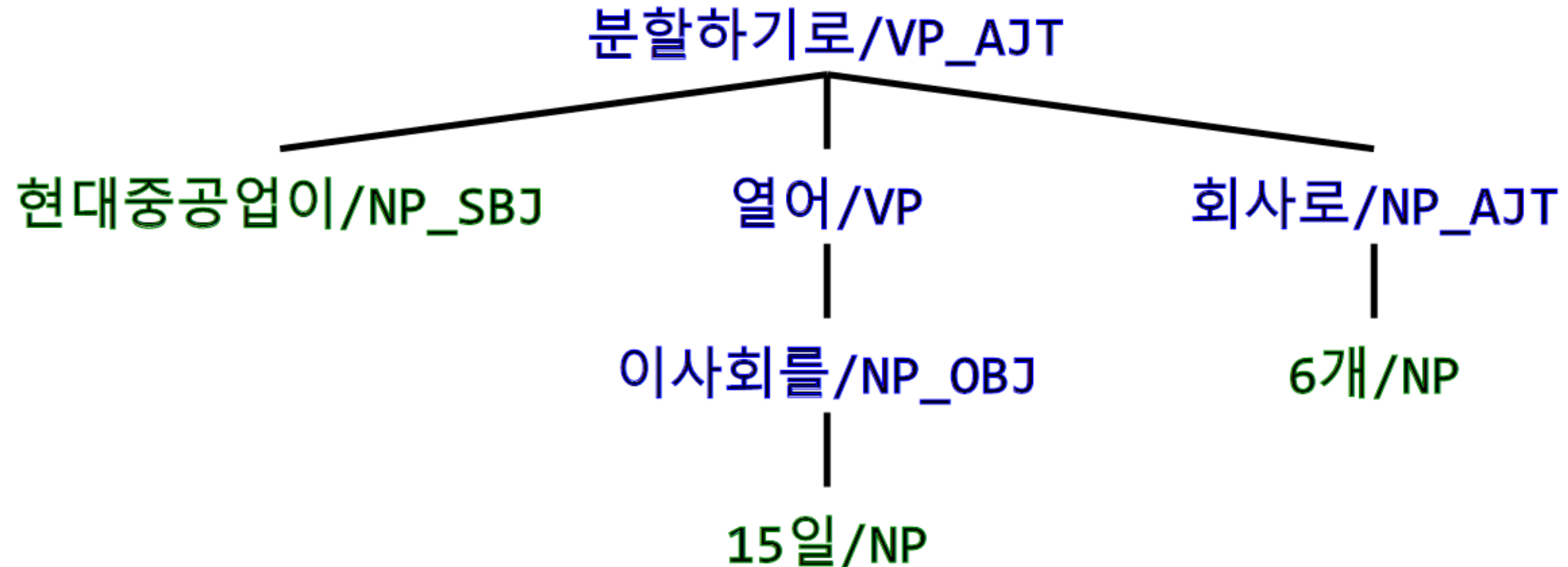
Componentized Output Tree



Tree Visualization (3/3)

49

Reassembled Output Tree



Serving Model

50

- TensorFlow Serving
 - Server-client based API model
 - Requires custom patch from git
 - Still buggy for SyntaxNet
- SyntaxNet demo scripts
 - Pipe from stdin/stdout
 - Messy, but most reliable

Example (1/5)

51

- Python-hosted script on Amazon EC2 which invokes SyntaxNet
- Generates diagram of dependency parsing tree from CoNLL-U output
- Modified tree generator to output SVG format for flexibility
- http://andrewmatteson.name/psg_tree.htm

Example (2/5)

52

Sample Sentence

- 현대중공업이 15일 이사회를 열어 6개 회사로 분할하기로 한 것은 극심한 실적 부진을 겪고 있는 조선·해양 부문에 회사의 자원과 역량을 집중시키고 비조선 부문은 자생력을 키워주기 위한 '고육지책'으로 풀이된다.

Example (3/5)

53

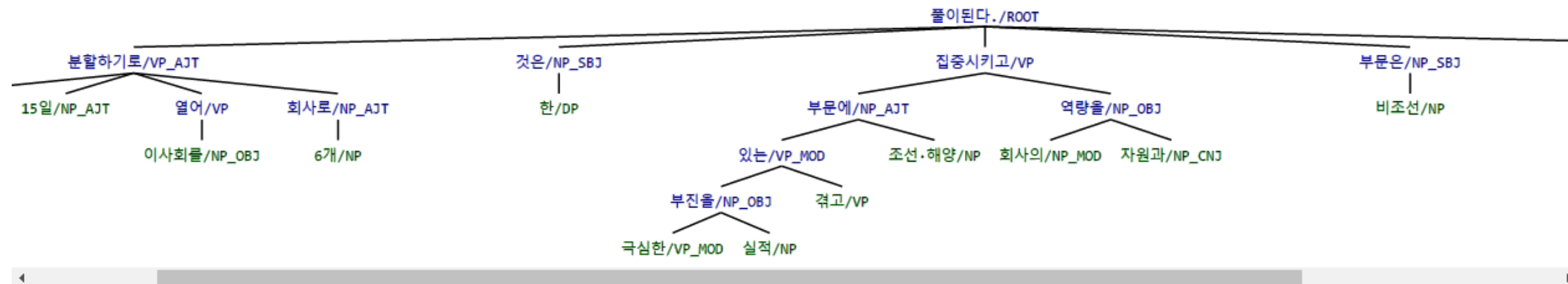
SyntaxNet Tree Generator (Korean/Komoran3)

Sentence to Parse:

현대중공업이 15일 이사회를 열어 6개 회사로 분할하기로 한 것은 극심한 실적 부진을 겪고 있는 조선·해양 부문에 회사의 자원과 역량을 집중시키고 비조선 부문은 자생력을 키워주기 위한 '고육지책'으로 풀이된다.

Submit

Visualization



Example (4/5)

54

CoNLL-U Output

Details

1 POS Tagging

2 Dependency Parsing

3 Reassembly

4 Semantic Analysis

D Debug Logs

ID	FORM	LEMMA	UPOSTAG	XPOSTAG	FEATS	HEAD	DEPREL
1	현대중공업이	현대중공업/NNP + 이/JKS	-	-	-	7	NP_SBJ
2	15일	15/SN + 일/NNB	-	-	-	7	NP_AJT
3	이사회를	이사회/NNG + 를/JKO	-	-	-	4	NP_OBJ
4	열어	열/VV + 어/EC	-	-	-	7	VP
5	6개	6/SN + 개/NNB	-	-	-	6	NP
6	회사로	회사/NNG + 로/JKB	-	-	-	7	NP_AJT
7	분할하기로	분할/NNG + 하/XSV + 기/ETN + 로/JKB	-	-	-	27	VP_AJT
8	한	한/MM	-	-	-	9	DP
9	것은	것/NNB + 은/JX	-	-	-	27	NP_SBJ
10	극심한	극심/XR + 하/XSA + ㄴ/ETM	-	-	-	12	VP_MOD
11	실적	실적/NNG	-	-	-	12	NP
12	부진을	부진/NNG + 을/JKO	-	-	-	14	NP_OBJ

Example (5/5)

55

1 POS Tagging

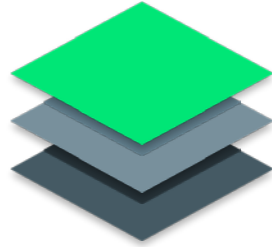
2 Dependency Parsing

3 Reassembly

4 Semantic Analysis

D Debug Logs

ID	FORM	DEPREL	Categories	Synonyms
1	현대중공업이	NP_SBJ	<div>1978년 해체</div> <div>중립성에 이의가 제기된 문서/2013년 5월</div> <div>대한민국의 조선업 기업</div> <div>대한민국의 중공업 기업</div> <div>한국거래소 상장 기업</div> <div>현대중공업그룹</div> <div>출처가 필요한 글/2016년 2월</div> <div>대한민국의 군수산업체</div> <div>1978년 설립</div> <div>1973년 설립</div>	Hyundai Heavy Industries 현대조선중공업 현대중공업
2	15일	NP_AJT	<div>등음이의어 문서</div>	
3	이사회를	NP_OBJ	<div>회사법</div> <div>상법</div> <div>경영</div>	Board of directors 이사회
4	열어	VP		



감사합니다.

Andrew Matteson
E-mail : amatteson@korea.ac.kr