

LING/C SC 581:

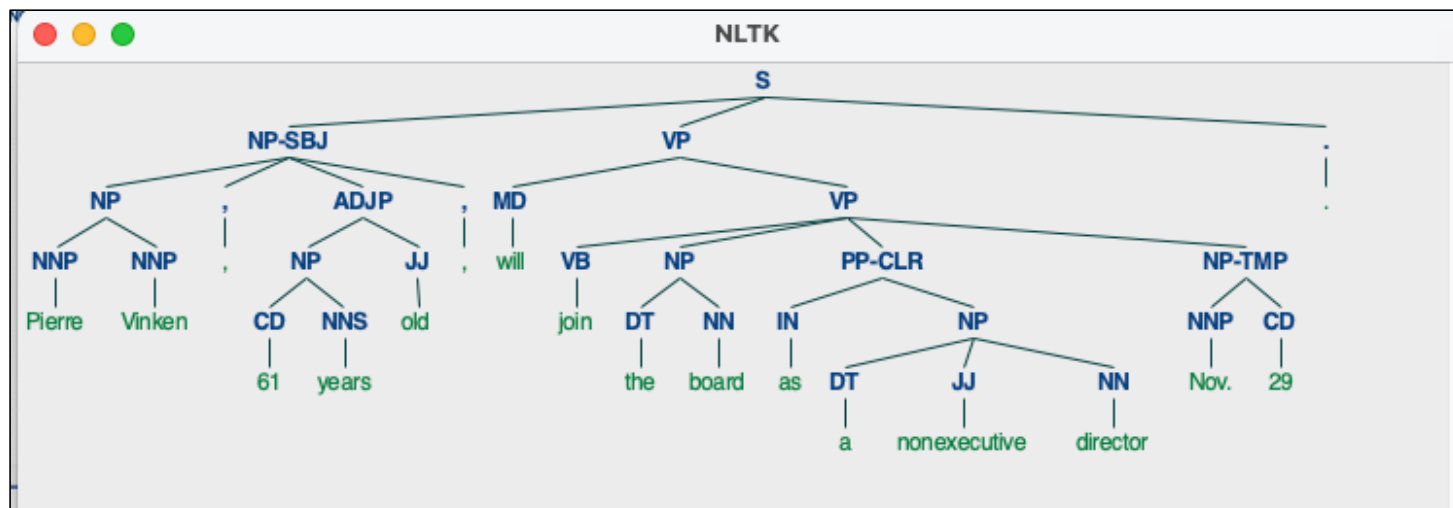
Advanced Computational Linguistics

Lecture 21

Today's Topic

- Copy the full Penn Treebank (PTB) corpus from the course website
 - instructions given out in Panopto and in class (**not on class slides!**)
- Homework 9: install tregex
 - we'll be using the software called tregex to search the treebanks
 - it's written in Java and requires a Java runtime environment

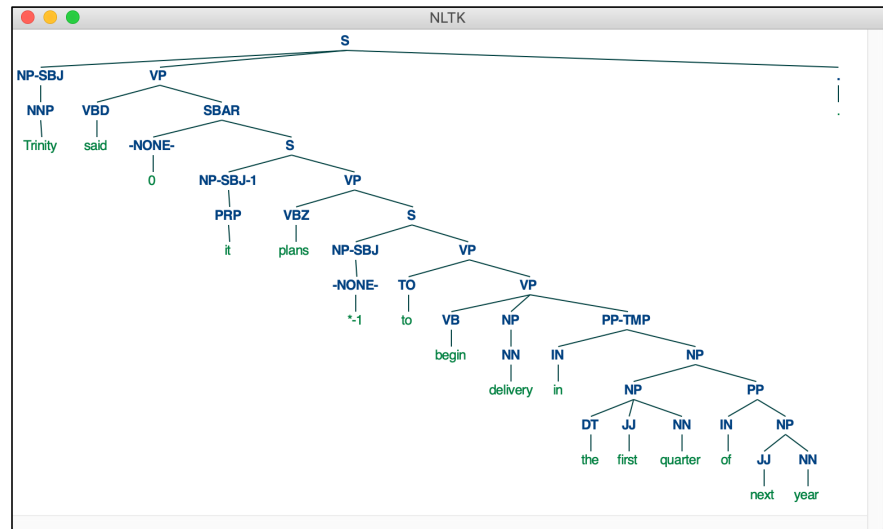
nltk



```
Tree('S', [Tree('NP-SBJ', [Tree('NP', [Tree('NNP', ['Pierre']), Tree('NNP', ['Vinken'])]), Tree(',', ['']), Tree('ADJP', [Tree('NP', [Tree('CD', ['61']), Tree('NNS', ['years'])]), Tree('JJ', ['old'])]), Tree(',', [''])]), Tree('VP', [Tree('MD', ['will']), Tree('VP', [Tree('VB', ['join']), Tree('NP', [Tree('DT', ['the']), Tree('NN', ['board'])]), Tree('PP-CLR', [Tree('IN', ['as']), Tree('NP', [Tree('DT', ['a']), Tree('JJ', ['nonexecutive']), Tree('NN', ['director'])])]), Tree('NP-TMP', [Tree('NNP', ['Nov.']), Tree('CD', ['29'])])])]), Tree('.', ['.'])])])])
```

nltk

- `>>> t[-1].draw()`
- The *sample* of the well-known Penn Treebank (PTB) Wall Street Journal (WSJ) corpus includes:
 - 3,914 parsed sentences
 - 49,000+ parsed sentences in the full corpus



nltk

- Words:

```
>>> w = treebank.words()
```

```
>>> len(w)
```

```
100676
```

```
>>> w
```

```
['Pierre', 'Vinken', ',', '61', 'years', 'old', ',', ...]
```

```
>>> tw = treebank.tagged_words()
```

```
>>> tw
```

```
[('Pierre', 'NNP'), ('Vinken', 'NNP'), (',', ',', ','), ...]
```

```
>>> tw[:10]
```

```
[('Pierre', 'NNP'), ('Vinken', 'NNP'), (',', ',', ','), ('61', 'CD'), ('years', 'NNS'), ('old', 'JJ'), (',', ',', ','), ('will', 'MD'), ('join', 'VB'), ('the', 'DT')]
```

nltk: Corpus Readers

- <http://www.nltk.org/howto/corpus.html#parsed-corpora>
 - The NLTK data package includes a 10% sample of the Penn Treebank (in `treebank`), as well as the Sinica Treebank (in `sinica_treebank`).
- Reading the Penn Treebank (Wall Street Journal sample):

```
>>> from nltk.corpus import treebank
1. treebank.fileids()
2. treebank.words(fileid)
3. treebank.tagged_words(fileid)
4. treebank.parsed_sents(fileid)
```

nltk: Corpus Readers

treebank.words(*fileid*)

```
>>> len(treebank.words('wsj_0003.mrg'))
782
>>> treebank.words('wsj_0003.mrg')
['A', 'form', 'of', 'asbestos', 'once', 'used', '*', ...]
>>> list(treebank.words('wsj_0003.mrg'))[:200]
['A', 'form', 'of', 'asbestos', 'once', 'used', '*', '*', 'to', 'make', 'Kent', 'cigarette', 'filters', 'has', 'caused',
',', 'a', 'high', 'percentage', 'of', 'cancer', 'deaths', 'among', 'a', 'group', 'of', 'workers', 'exposed', '*', 'to',
'it', 'more', 'than', '30', 'years', 'ago', ',', 'researchers', 'reported', '0', '*T*-1', '.', 'The', 'asbestos', 'fib
er', ',', 'crocidolite', ',', 'is', 'unusually', 'resilient', 'once', 'it', 'enters', 'the', 'lungs', ',', 'with', 'ev
en', 'brief', 'exposures', 'to', 'it', 'causing', 'symptoms', 'that', '*T*-1', 'show', 'up', 'decades', 'later', ',',
'researchers', 'said', '0', '*T*-2', '.', 'Lorillard', 'Inc.', ',', 'the', 'unit', 'of', 'New', 'York-based', 'Loews',
'Corp.', 'that', '*T*-2', 'makes', 'Kent', 'cigarettes', ',', 'stopped', 'using', 'crocidolite', 'in', 'its', 'Micron
ite', 'cigarette', 'filters', 'in', '1956', '.', 'Although', 'preliminary', 'findings', 'were', 'reported', '*-2', 'mo
re', 'than', 'a', 'year', 'ago', ',', 'the', 'latest', 'results', 'appear', 'in', 'today', "'s", 'New', 'England', 'Jo
urnal', 'of', 'Medicine', ',', 'a', 'forum', 'likely', '*', 'to', 'bring', 'new', 'attention', 'to', 'the', 'problem',
',', 'A', 'Lorillard', 'spokewoman', 'said', ',', '``', 'This', 'is', 'an', 'old', 'story', '.', 'We', "'re", 'talkin
g', 'about', 'years', 'ago', 'before', 'anyone', 'heard', 'of', 'asbestos', 'having', 'any', 'questionable', 'properti
es', '.', 'There', 'is', 'no', 'asbestos', 'in', 'our', 'products', 'now', '.', "'", 'Neither', 'Lorillard', 'nor', '
the', 'researchers', 'who', '*T*-3', 'studied', 'the', 'workers', 'were', 'aware', 'of', 'any', 'research', 'on', 'smo
kers', 'of', 'the', 'Kent', 'cigarettes', '.']
>>> █
```


nlTK: Corpus Readers

treebank.tagged_words(*fileid*)

```
>>> >>> list(treebank.tagged_words('wsj_0003.mrg'))[:100]
[('A', 'DT'), ('form', 'NN'), ('of', 'IN'), ('asbestos', 'NN'), ('once', 'RB'), ('used', 'VBN'), (*, '-NONE-'), (*, '-NONE-'), ('to', 'TO'), ('make', 'VB'), ('Kent', 'NNP'), ('cigarette', 'NN'), ('filters', 'NNS'), ('has', 'VBZ'), ('caused', 'VBN'), ('a', 'DT'), ('high', 'JJ'), ('percentage', 'NN'), ('of', 'IN'), ('cancer', 'NN'), ('deaths', 'NNS'), ('among', 'IN'), ('a', 'DT'), ('group', 'NN'), ('of', 'IN'), ('workers', 'NNS'), ('exposed', 'VBN'), (*, '-NONE-'), ('to', 'TO'), ('it', 'PRP'), ('more', 'RBR'), ('than', 'IN'), ('30', 'CD'), ('years', 'NNS'), ('ago', 'IN'), (',', ','), ('researchers', 'NNS'), ('reported', 'VBD'), ('0', '-NONE-'), (*T*-1, '-NONE-'), ('.', '.'), ('The', 'DT'), ('asbestos', 'NN'), ('fiber', 'NN'), (',', ','), ('crocidolite', 'NN'), (',', ','), ('is', 'VBZ'), ('unusually', 'RB'), ('resilient', 'JJ'), ('once', 'IN'), ('it', 'PRP'), ('enters', 'VBZ'), ('the', 'DT'), ('lungs', 'NNS'), (',', ','), ('with', 'IN'), ('even', 'RB'), ('brief', 'JJ'), ('exposures', 'NNS'), ('to', 'TO'), ('it', 'PRP'), ('causing', 'VBG'), ('symptoms', 'NNS'), ('that', 'WDT'), (*T*-1, '-NONE-'), ('show', 'VBP'), ('up', 'RP'), ('decades', 'NNS'), ('later', 'JJ'), (',', ','), ('researchers', 'NNS'), ('said', 'VBD'), ('0', '-NONE-'), (*T*-2, '-NONE-'), ('.', '.'), ('Lorillard', 'NNP'), ('Inc.', 'NNP'), (',', ','), ('the', 'DT'), ('unit', 'NN'), ('of', 'IN'), ('New', 'JJ'), ('York-based', 'JJ'), ('Loews', 'NNP'), ('Corp.', 'NNP'), ('that', 'WDT'), (*T*-2, '-NONE-'), ('makes', 'VBZ'), ('Kent', 'NNP'), ('cigarettes', 'NNS'), (',', ','), ('stopped', 'VBD'), ('using', 'VBG'), ('crocidolite', 'NN'), ('in', 'IN'), ('its', 'PRP$'), ('Micronite', 'NN'), ('cigarette', 'NN'), ('filters', 'NNS')]
>>> █
```


Penn POS Tagset

POS Tag	Description	Example
CC	coordinating conjunction	and
CD	cardinal number	1, third
DT	determiner	the
EX	existential there	<i>there</i> is
FW	foreign word	d'hoevre
IN	preposition/subordinating conjunction	in, of, like
JJ	adjective	green
JJR	adjective, comparative	greener
JJS	adjective, superlative	greenest
LS	list marker	1)
MD	modal	could, will
NN	noun, singular or mass	table
NNS	noun plural	tables
NNP	proper noun, singular	John
NNPS	proper noun, plural	Vikings

Penn POS Tagset

PDT	predeterminer	<i>both the boys</i>
POS	possessive ending	<i>friend's</i>
PRP	personal pronoun	<i>I, he, it</i>
PRPS	possessive pronoun	<i>my, his</i>
RB	adverb	<i>however, usually, naturally, here, good</i>
RBR	adverb, comparative	<i>better</i>
RBS	adverb, superlative	<i>best</i>
RP	particle	<i>give up</i>
TO	to	<i>to go, to him</i>
UH	interjection	<i>uhhuhhuhh</i>
VB	verb, base form	<i>take</i>
VBD	verb, past tense	<i>took</i>
VBG	verb, gerund/present participle	<i>taking</i>
VBN	verb, past participle	<i>taken</i>
VBP	verb, sing. present, non-3d	<i>take</i>
VBZ	verb, 3rd person sing. present	<i>takes</i>

Penn POS Tagset

WDT	wh-determiner	which
WP	wh-pronoun	who, what
WP\$	possessive wh-pronoun	whose
WRB	wh-abverb	where, when

Penn Syntax Tagset

- (from The Penn Treebank: An overview, Taylor, Marcus & Santorini)

Table 1.2. The Penn Treebank syntactic tagset

ADJP	Adjective phrase
ADVP	Adverb phrase
NP	Noun phrase
PP	Prepositional phrase
S	Simple declarative clause
SBAR	Subordinate clause
SBARQ	Direct question introduced by <i>wh</i> -element
SINV	Declarative sentence with subject-aux inversion
SQ	Yes/no questions and subconstituent of SBARQ excluding <i>wh</i> -element
VP	Verb phrase
WHADVP	Wh-adverb phrase
WHNP	Wh-noun phrase
WHPP	Wh-prepositional phrase
X	Constituent of unknown or uncertain category
*	“Understood” subject of infinitive or imperative
0	Zero variant of <i>that</i> in subordinate clauses
T	Trace of <i>wh</i> -Constituent

Penn Syntax Tagset

- (from The Penn Treebank: An overview, Taylor, Marcus & Santorini)

<i>Text Categories</i>	
-HLN	headlines and datelines
-LST	list markers
-TTL	titles
<i>Grammatical Functions</i>	
-CLF	true clefts
-NOM	non NPs that function as NPs
-ADV	clausal and NP adverbials
-LGS	logical subjects in passives
-PRD	non VP predicates
-SBJ	surface subject
-TPC	topicalized and fronted constituents
-CLR	closely related - see text
<i>Semantic Roles</i>	
-VOC	vocatives
-DIR	direction & trajectory
-LOC	location
-MNR	manner
-PRP	purpose and reason
-TMP	temporal phrases

Tregex

- URL: <https://nlp.stanford.edu/software/tregex.shtml>

Contents

The download is a 9 Mb zip file. It contains:

1. README-tregex.txt -- Basic information about the distribution, including a "quickstart" guide.
2. README-tsurgeon.txt -- information about Tsurgeon.
3. README-gui.txt -- information about using the graphical interface
4. LICENSE -- Tregex is licensed under the Gnu General Public License.
5. stanford-tregex.jar -- This is a JAR file containing all the Stanford classes necessary to run tregex.
6. src directory -- a directory with the source files for Tregex and Tsurgeon
7. lib directory -- library files required for recompiling the distribution (with Mac OS X customization; see [lib/ABOUT-AppleJavaExtensions.txt](#) for removing this dependency)
8. build.xml, Makefile -- files for recompiling (with ant or make) the distribution
9. javadoc -- Javadocs for the distribution
10. tregex.sh, tsurgeon.sh -- sample scripts for running Tregex and Tsurgeon from the command line
11. run-tregex-gui.command, run-tregex-gui.bat -- shell script for running the graphical interface for Tregex with more memory for searching larger treebanks; can be double-clicked to open on a Mac or PC, respectively
12. examples directory -- example files for Tregex and Tsurgeon

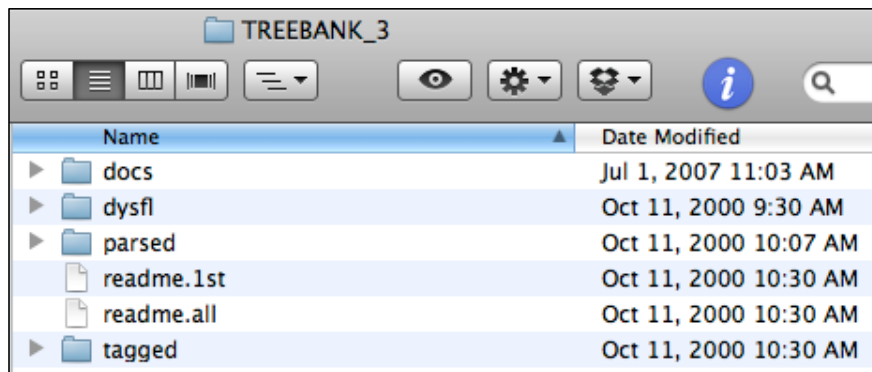
Download

[Download Tregex version 4.2.0](#) (source and executables for all platforms)

[Download Tregex version 3.4 Mac OS X disk image](#) (GUI packaged as Mac application; Java 1.7 runtime included)

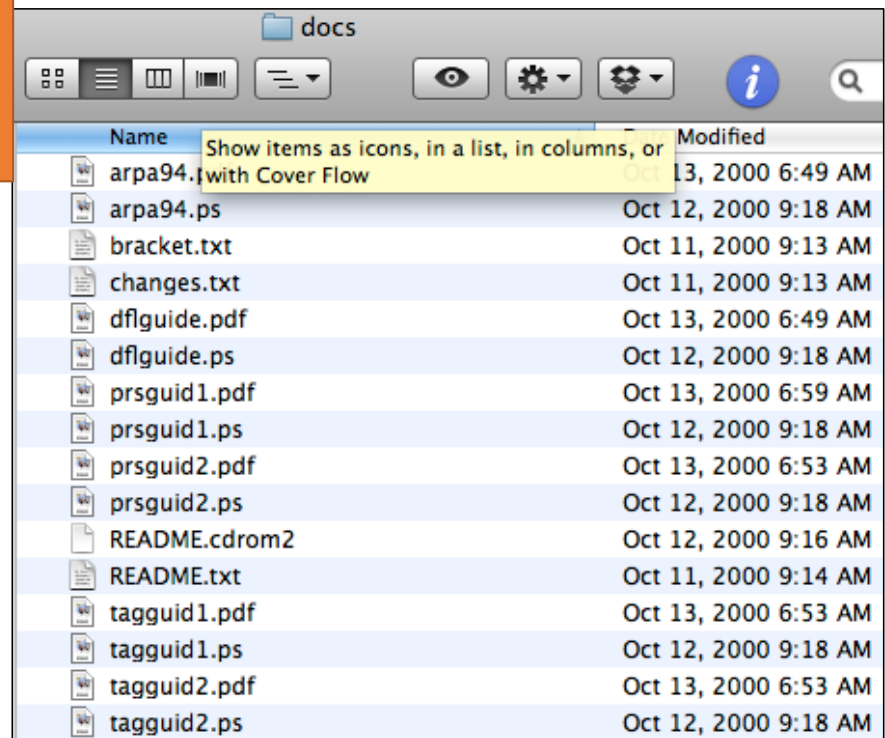
Trebank Guides

- Tagging Guide
- arpa94 paper
- Parse Guide



File browser window showing the directory structure of TREEBANK_3. The window title is "TREEBANK_3". The table below lists the contents of the directory.

Name	Date Modified
▶ docs	Jul 1, 2007 11:03 AM
▶ dysfl	Oct 11, 2000 9:30 AM
▶ parsed	Oct 11, 2000 10:07 AM
readme.1st	Oct 11, 2000 10:30 AM
readme.all	Oct 11, 2000 10:30 AM
▶ tagged	Oct 11, 2000 10:30 AM



File browser window showing the contents of the docs directory. The window title is "docs". A tooltip is visible over the first row: "Show items as icons, in a list, in columns, or". The table below lists the contents of the directory.

Name	Modified
arpa94.twith Cover Flow	Oct 13, 2000 6:49 AM
arpa94.ps	Oct 12, 2000 9:18 AM
bracket.txt	Oct 11, 2000 9:13 AM
changes.txt	Oct 11, 2000 9:13 AM
dflguide.pdf	Oct 13, 2000 6:49 AM
dflguide.ps	Oct 12, 2000 9:18 AM
prsguid1.pdf	Oct 13, 2000 6:59 AM
prsguid1.ps	Oct 12, 2000 9:18 AM
prsguid2.pdf	Oct 13, 2000 6:53 AM
prsguid2.ps	Oct 12, 2000 9:18 AM
README.cdrom2	Oct 12, 2000 9:16 AM
README.txt	Oct 11, 2000 9:14 AM
tagguid1.pdf	Oct 13, 2000 6:53 AM
tagguid1.ps	Oct 12, 2000 9:18 AM
tagguid2.pdf	Oct 13, 2000 6:53 AM
tagguid2.ps	Oct 12, 2000 9:18 AM

Treebank Guides

- Parts-of-speech (POS) Tagging Guide, tagguid1.pdf (34 pages):

Part-of-Speech Tagging Guidelines
for the Penn Treebank Project
(3rd Revision, 2nd printing)

Beatrice Santorini

June 1990 ¹

tagguid2.pdf: addendum, see POS tag 'TO'

Treebank Guides

- Parsing guide 1, prsguid1.pdf (318 pages):

**Bracketing Guidelines for Treebank II Style
Penn Treebank Project ¹**

Principal authors:

Ann Bies, Mark Ferguson, Karen Katz, and Robert MacIntyre

Major contributors:

Victoria Tredinnick, Grace Kim, Mary Ann Marcinkiewicz, Britta Schasberger ²

January 1995

prsguid2.pdf: addendum for the Switchboard corpus

tregex

1. Shell file:

```
run-tregex-gui.command
#!/bin/sh
java -mx300m -cp `dirname $0`/stanford-tregex.jar edu.stanford.nlp.trees.tregex.gui.TregexGUI
```

2. Select the PTB directory

- TREEBANK_3/parsed/mrg/ws_j/
- *you can select more directories*

3. Browse Trees

The image illustrates the Tregex GUI workflow. It starts with a terminal window showing the command to run the Tregex GUI. The main window shows a file browser on the left with a directory tree where 'wsj' is selected. The central part of the window is the search interface, including a search pattern field, a 'Search' button, and a 'Browse Trees' button. On the right, a list of matches is displayed, with the first match highlighted: 'wsj_0001.mrg-1 Pierre Vinken, 61 years old, will join the board'. Below the search interface, a detailed parse tree is shown for the sentence 'Vinken, 61 years old, will join the board as a nonexecutive director Nov. 29.' The tree structure includes nodes like NP-SBJ, VP, MD, and PP-CLR, with the corresponding words and punctuation below them.

Possible macOS Problem

- Disk image version, the Java runtime environment seems to pick wrong fonts. Display is hard to read.

The screenshot shows a Java application window titled "Tregex". The window is divided into several sections:

- Search pattern:** Contains a text field with "Rece" and a dropdown menu.
- Matches:** A list of search results. The fifth result, "ws_0347mg-5The sea bone me e h g s w o n d i t", is highlighted in blue.
- Parse Tree:** A complex tree diagram showing the syntactic structure of the sentence "The sea bone me e h g s w o n d i t". The root node is "S", which branches into "NP6BJ" and "VP". The "NP6BJ" node branches into "DT" (The) and "JJ" (sea). The "VP" node branches into "NNS" (bone) and "MDRB" (me). The "NNS" node branches into "VB" (e) and "NP". The "MDRB" node branches into "VP" (h) and "NP" (g). The "NP" node branches into "PP" (s) and "NP". The "PP" node branches into "P" (w) and "NP". The "NP" node branches into "NNSN" (o) and "NP". The "NNSN" node branches into "p" and "n". The "NP" node branches into "NP" (n) and "NP" (d). The "NP" node branches into "DT" (a) and "JJ" (t). The "NP" node branches into "NN" (s) and "NP". The "DT" node branches into "a". The "JJ" node branches into "b". The "NN" node branches into "m". The "NP" node branches into "JJ" (e) and "NN" (x). The "JJ" node branches into "e". The "NN" node branches into "p". The "NP" node branches into "CC" (h) and "NN" (g). The "CC" node branches into "h". The "NN" node branches into "o". The "NP" node branches into "NN" (s) and "NP". The "NN" node branches into "u". The "NP" node branches into "WHFP1" (n) and "NP". The "WHFP1" node branches into "at". The "NP" node branches into "WHNP" (t) and "NP". The "WHNP" node branches into "WDT" (d) and "DT" (t). The "WHNP" node branches into "t". The "NP" node branches into "NP6BJ2" (n) and "NP". The "NP6BJ2" node branches into "NN" (n) and "NP". The "NP6BJ2" node branches into "n". The "NP" node branches into "VP" (z) and "NP". The "VP" node branches into "VBZ" (b) and "NP". The "VP" node branches into "ADJPPRD" (y) and "NP". The "ADJPPRD" node branches into "JJ" (o) and "NP". The "ADJPPRD" node branches into "n". The "NP" node branches into "S" (n) and "NP". The "S" node branches into "JJ" (n) and "NP". The "S" node branches into "key". The "NP" node branches into "NP6BJ8" (n) and "NP". The "NP6BJ8" node branches into "NP" (n) and "NP". The "NP6BJ8" node branches into "n". The "NP" node branches into "VP" (n) and "NP". The "VP" node branches into "VP" (n) and "NP". The "VP" node branches into "n". The "NP" node branches into "NP" (n) and "NP". The "NP" node branches into "n". The "NP" node branches into "NP" (n) and "NP". The "NP" node branches into "n".

Possible macOS Problem

- If this happens, download the non-image link

[Download Tregex version 4.2.0](#) (source and executables for all platforms)

[Download Tregex version 3.4 Mac OS X disk image](#) (GUI packaged as Mac application; Java 1.7 runtime included)

Name	Date Modified	Size	Kind
build.xml	Nov 17, 2020 at 2:57 AM	6 KB	XML Document
examples	Nov 17, 2020 at 2:57 AM	--	Folder
lib	Nov 17, 2020 at 2:57 AM	--	Folder
LICENSE.txt	Nov 17, 2020 at 2:57 AM	18 KB	Plain Text
Makefile	Nov 17, 2020 at 2:57 AM	567 bytes	Makefile
README-gui.txt	Nov 17, 2020 at 2:57 AM	11 KB	Plain Text
README-tregex.txt	Nov 17, 2020 at 2:57 AM	16 KB	Plain Text
README-tsurgeon.txt	Nov 17, 2020 at 2:57 AM	18 KB	Plain Text
run-tregex-gui.bat	Nov 17, 2020 at 2:57 AM	83 bytes	Document
run-tregex-gui.command	Nov 17, 2020 at 2:57 AM	104 bytes	Termina...ll script
Semgrex.ppt	Nov 17, 2020 at 2:57 AM	386 KB	PowerP...n (.ppt)
stanford-tregex-4.2.0-javadoc.jar	Nov 17, 2020 at 2:57 AM	3.4 MB	Java JAR file
stanford-tregex-4.2.0-sources.jar	Nov 17, 2020 at 2:57 AM	2.6 MB	Java JAR file
stanford-tregex-4.2.0.jar	Nov 17, 2020 at 2:57 AM	2.9 MB	Java JAR file
stanford-tregex.jar	Nov 17, 2020 at 2:57 AM	2.9 MB	Java JAR file
tregex.sh	Nov 17, 2020 at 2:57 AM	134 bytes	Shell Script

for macOS, run this command

tregex

- Search

- NP-SBJ << (*dominates*) vs. < (*immediately dominates*) NNP

Pattern	Trees Matched	Total Matches
NP-SBJ << NNP	19862	53523
NP-SBJ < NNP	11994	22740

Search pattern: NP-SBJ << NNP
Recent: NP-SBJ << NNP
Pattern: NP-SBJ << NNP
Tree size: [slider]
Tsurgeon script:
Matches: wsj_0001.mrg-1 Pierre. ...
Match stats: 19862 unique trees found with 53523 total matches.

EBANK_3/parsed/mrg/ws/00/wsj_0001.mrg

```
graph TD
    S --- NP_SBJ[NP-SBJ]
    S --- VP1[VP]
    NP_SBJ --- NP1[NP]
    NP_SBJ --- ADJP[ADJP]
    NP1 --- NNP1[NNP]
    NP1 --- NNP2[NNP]
    ADJP --- CD[CD]
    ADJP --- NNS[NNS]
    ADJP --- JJ[JJ]
    ADJP --- old[old]
    NP1 --- Pierre[Pierre]
    NP1 --- Vinken[Vinken]
    CD --- 61[61]
    NNS --- years[years]
    VP1 --- MD[MD]
    VP1 --- VB[VB]
    MD --- will[will]
    VP1 --- VP2[VP]
    VP2 --- join[join]
    VP2 --- DT1[DT]
    VP2 --- NN1[NN]
    DT1 --- the[the]
    NN1 --- board[board]
    VP2 --- IN[IN]
    IN --- as[as]
    VP2 --- PP_CLR[PP-CLR]
    PP_CLR --- DT2[DT]
    PP_CLR --- NP2[NP]
    DT2 --- a[a]
    NP2 --- JJ1[JJ]
    NP2 --- NN2[NN]
    JJ1 --- nonexecutive[nonexecutive]
    NP2 --- director[director]
    VP2 --- NP_TMP[NP-TMP]
    NP_TMP --- NNP3[NNP]
    NP_TMP --- CD2[CD]
    NNP3 --- Nov[Nov]
    CD2 --- 29[29]
```

Pierre Vinken , 61 years old , will join the board as a nonexecutive director Nov. 29 .

Pattern: NP-SBJ < NNP
Tree size: [slider]
Tsurgeon script:
Match stats: 11994 unique trees found with 22740 total matches.

BANK_3/parsed/mrg/ws/00/wsj_0001.mrg

```
graph TD
    S --- NP_SBJ[NP-SBJ]
    S --- VP[VP]
    NP_SBJ --- NNP1[NNP]
    NP_SBJ --- NNP2[NNP]
    NNP1 --- Mr[Mr.]
    NNP2 --- Vinken[Vinken]
    VP --- VBZ[VBZ]
    VBZ --- is[is]
    VP --- NP_PRD[NP-PRD]
    NP_PRD --- NP3[NP]
    NP_PRD --- PP[PP]
    NP3 --- NN[NN]
    NN --- chairman[chairman]
    PP --- of[of]
    PP --- NP4[NP]
    NP4 --- NP5[NP]
    NP4 --- DT[DT]
    NP5 --- NNP4[NNP]
    NP5 --- NNP5[NNP]
    DT --- the[the]
    NNP4 --- Elsevier[Elsevier]
    NNP5 --- N.V.[N.V.]
    NP4 --- NN6[NN]
    NN6 --- Dutch[Dutch]
    NP4 --- VBG[VBG]
    VBG --- publishing[publishing]
    NP4 --- NN7[NN]
    NN7 --- group[group]
```

Mr. Vinken is chairman of Elsevier N.V. , the Dutch publishing group .

tregex

- README-tregex.txt

Tregex Pattern Syntax and Uses

Using a Tregex pattern, you can find only those trees that match the pattern you're looking for. The following table shows the symbols that are allowed in the pattern, and below there is more information about using these patterns.

Symbol	Meaning
A << B	A dominates B
A >> B	A is dominated by B
A < B	A immediately dominates B
A > B	A is immediately dominated by B
A \$ B	A is a sister of B (and not equal to B)
A .. B	A precedes B
A . B	A immediately precedes B
A ,, B	A follows B
A , B	A immediately follows B
A <<, B	B is a leftmost descendent of A
A <<- B	B is a rightmost descendent of A
A >>, B	A is a leftmost descendent of B
A >>- B	A is a rightmost descendent of B
A <, B	B is the first child of A
A >, B	A is the first child of B
A <- B	B is the last child of A
A >- B	A is the last child of B
A <# B	B is the last child of A
A ># B	A is the last child of B
A <i B	B is the ith child of A (i > 0)
A >i B	A is the ith child of B (i > 0)
A <-i B	B is the ith-to-last child of A (i > 0)
A >-i B	A is the ith-to-last child of B (i > 0)

A <: B	B is the only child of A
A >: B	A is the only child of B
A <<: B	A dominates B via an unbroken chain (length > 0) of unary local trees.
A >>: B	A is dominated by B via an unbroken chain (length > 0) of unary local trees.
A \$++ B	A is a left sister of B (same as \$.. for context-free trees)
A \$-- B	A is a right sister of B (same as \$., for context-free trees)
A \$+ B	A is the immediate left sister of B (same as \$. for context-free trees)
A \$- B	A is the immediate right sister of B (same as \$, for context-free trees)
A \$.. B	A is a sister of B and precedes B
A \$., B	A is a sister of B and follows B
A \$. B	A is a sister of B and immediately precedes B
A \$, B	A is a sister of B and immediately follows B
A <+(C) B	A dominates B via an unbroken chain of (zero or more) nodes matching description C
A >+(C) B	A is dominated by B via an unbroken chain of (zero or more) nodes matching description C
A .+(C) B	A precedes B via an unbroken chain of (zero or more) nodes matching description C
A ,+(C) B	A follows B via an unbroken chain of (zero or more) nodes matching description C
A <<# B	B is a head of phrase A
A >># B	A is a head of phrase B
A <# B	B is the immediate head of phrase A
A ># B	A is the immediate head of phrase B
A == B	A and B are the same node
A : B	[this is a pattern-segmenting operator that places no constraints on the relationship between A and B]