

Gán nhãn từ loại (Part-of-speech tagging)

Read Chapter 8 - Speech and
Language Processing

1

Definition

- Part of Speech (POS) tagging: assign each word in a sentence with an appropriate POS.
 - Input: a string of words + a tagset
 - Output: a best tag for each word

[Example 1](#)
[Example 2](#)
[Example 3](#)
[Example 4](#)
[Example 5](#)

- Tagging makes parsing easier

2

Why POS tagging?

- **Simple:** can be done by many different methods
 - Can be done well with methods that look at local context
 - Though should “really” do it by parsing!
- **Applications:**
 - Text-to-speech: **record** - N: [ˈreko:d], V: [riˈko:d]; **lead** – N [led], V: [li:d]
 - Can be a preprocessor for a parser. The parser can do it better but more expensive
 - Speech recognition, parsing, information retrieval, etc.
- **Easy to evaluate** (*how many tags are correct?*)

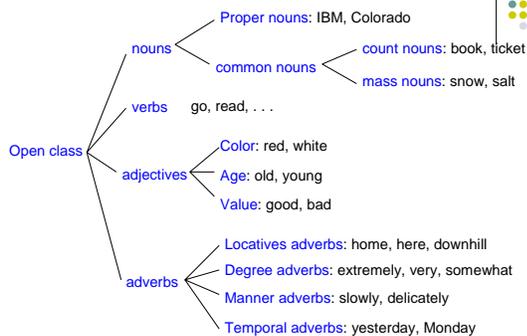
3

English word classes

- **Closed class** (function words): fixed membership
 - Prepositions: on, under, over,...
 - Particles: abroad, about, around, before, in, instead, since, without,...
 - Articles: a, an, the
 - Conjunctions: and, or, but, that,...
 - Pronouns: you, me, I, your, what, who,...
 - Auxiliary verbs: can, will, may, should,...
- **Open class:** new words can be added

4

English word classes



5

Tagsets for English

- 87 tags - Brown corpus
- Three most commonly used:
 - Small: 45 Tags - Penn treebank (next slide)
 - Medium size: 61 tags, British national corpus
 - Large: 146 tags, C7

6

Penn Treebank tags

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &</i>
CD	Cardinal number	<i>one, two, three</i>	TO	"to"	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential 'there'	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VBN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>wildest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WPS	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	"	Left quote	<i>(' or ")</i>
POS	Possessive ending	<i>'s</i>	"	Right quote	<i>(' or ")</i>
PP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>(, (, {, <</i>
PPS	Possessive pronoun	<i>your, one's</i>)	Right parenthesis	<i>(,), }, ></i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>(, ! ?)</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>(; ; ... - -)</i>
RP	Particle	<i>up, off</i>			

Example from Penn Treebank

- The grand jury commented on a number of other topics.
- ⇒ The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.

Problem of POS tagging

Problem of POS tagging is to resolve ambiguities, choosing the proper tag for the context.

Main types of taggers

- Stochastic tagging:** Maximum likelihood, Hidden Markov model tagging
Pr (Det-N) > Pr (Det-Det)
- Rule based tagging**
If <some pattern>
Then ... <some part of speech>

Approaches to Tagging

- HMM tagging:** 'Use all the information you have and guess'
- Constrain Grammar (CG) tagging:** 'Don't guess, just eliminate the impossible!'
- Transformation-based (TB) tagging:** 'Guess first, then change your mind if necessary!'

Stochastic POS tagging

For a given sentence or word sequence, pick the most likely tag for each word.

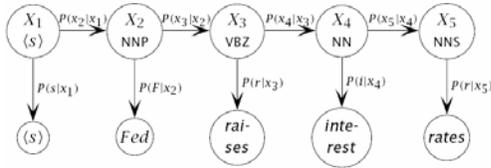
How?

- A Hidden Markov model (HMM) tagger:
Choose the tag sequence that maximizes:
 $P(\text{word}|\text{tag}) \cdot P(\text{tag}|\text{previous } n \text{ tags})$

The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.

$$\Rightarrow P(\text{jury}|\text{NN}) = 1/2$$

HMMs – POS example



Do supervised training, and then inference to decide POS tags (Bayesian network style)

13

HMM tagging

- **Bigram HMM Equation:** choose t_i for w_i that is most probably given t_{i-1} and w_i :

$$t_i = \operatorname{argmax}_j P(t_j | t_{i-1}, w_i) \quad (1)$$

- **A HMM simplifying assumption:** the tagging problem can be solved by looking at nearby words and tags.

$$t_i = \operatorname{argmax}_j P(t_j | t_{i-1}) P(w_i | t_j) \quad (2)$$

pr tag sequence \uparrow word (lexical) likelihood \uparrow
(tag co-occurrence)

14

Example

1. Secretariat/**NNP** is/**VBZ** expected/**VBN**
to/**TO** **race**/**VB** tomorrow/**NN**
2. People/**NNS** continue/**VBP** to/**TO** inquire/**VB**
the/**DT** reason/**NN** for/**IN** the/**DT** **race**/**NN**
for/**IN** outer/**JJ** space/**NN**

15

Suppose we have tagged all but **race**

- Look at just preceding word (bigram):
to/**TO** **race**/??? **NN** or **VB**?
the/**DT** **race**/???
- Applying (2): $t_i = \operatorname{argmax}_j P(t_j | t_{i-1}) P(w_i | t_j)$
- Choose tag with greater of the two probabilities:
 $P(\text{VB}|\text{TO})P(\text{race}|\text{VB})$ or $P(\text{NN}|\text{TO})P(\text{race}|\text{NN})$

16

Calculate Probabilities

Let's consider $P(\text{VB}|\text{TO})$ and $P(\text{NN}|\text{TO})$

- Can find these pr estimates by *counting* in a corpus (and normalizing)
- Expect that a **verb** is more likely to follow **TO** than a **Noun** is, since infinitives are common in English (*to race, to walk*). A **noun** can follow **TO** (*run to school*)
- From the Brown corpus
 $P(\text{NN}|\text{TO}) = .021$
 $P(\text{VB}|\text{TO}) = .340$

17

Calculate Probabilities

- $P(\text{race}|\text{VB})$ and $P(\text{race}|\text{NN})$: the *lexical likelihood* of the noun *race* given each tag, $P(\text{race}|\text{VB})$ and $P(\text{race}|\text{NN})$, e.g., "if we were expecting a verb, would it be *race*?"
- From the Brown corpus
 $P(\text{race}|\text{NN}) = 0.00041$
 $P(\text{race}|\text{VB}) = 0.00003$
- $P(\text{VB}|\text{TO})P(\text{race}|\text{VB}) = 0.00001$
 $P(\text{NN}|\text{TO})P(\text{race}|\text{NN}) = 0.000007$
- **race should be a VB after "TO"**

18

The full model

- Now we want the best sequence of tags for the whole sentence
- Given the sequence of words, W , we want to compute the most probably tag sequence, $T = t_1, t_2, \dots, t_n$ or,

$$\begin{aligned} \hat{T} &= \arg \max_{T \in \tau} P(T | W) \\ &= \arg \max_{T \in \tau} \frac{P(T)P(W | T)}{P(W)} \quad (\text{Bayes' Theorem}) \\ &= \arg \max_{T \in \tau} P(T)P(W | T) \end{aligned}$$

19

Expand this using chain rule

From chain rule for probabilities:

$$\begin{aligned} P(A, B) &= P(A|B)P(B) = P(B|A)P(A) \\ P(A, B, C) &= P(B, C|A)P(A) = P(C|A, B)P(B|A)P(A) \\ &= P(A)P(B|A)P(C|A, B) \\ P(A, B, C, D \dots) &= P(A)P(B|A)P(C|A, B)P(D|A, B, C \dots) \\ P(T)P(W | T) &= \prod_{i=1}^n P(w_i | w_1 t_1 \dots w_{i-1} t_{i-1}) P(t_i | w_1 t_1 \dots w_{i-1} t_{i-1}) \end{aligned}$$

← pr word
← tag history

20

Trigram assumption

- Probability of a word depends only on its tag

$$P(w_i | w_1 t_1 \dots t_{i-1} t_i) = P(w_i | t_i)$$

- Tag history approximated by two most recent tags (trigram: two most recent + current state)

$$P(t_i | w_1 t_1 \dots t_{i-1}) = P(t_i | t_{i-2} t_{i-1})$$

21

Replacing to the equation

$$P(T)P(W|T) =$$

$$P(t_1)P(t_2 | t_1) \prod_{i=3}^n P(t_i | t_{i-2} t_{i-1}) \left[\prod_{i=1}^n P(w_i | t_i) \right]$$

22

Estimate Probabilities

- Use relative frequencies from corpus to estimate these probabilities:

$$P(t_i | t_{i-1} t_{i-2}) = \frac{c(t_{i-2} t_{i-1} t_i)}{c(t_{i-2} t_{i-1})}$$

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

23

Problem

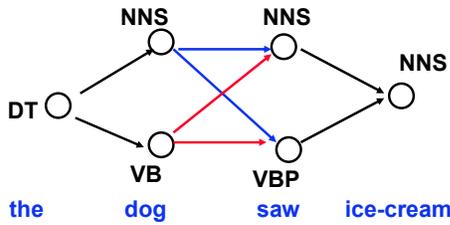
The problem to solve:

$$\hat{T} = \arg \max_{T \in \tau} P(T)P(W | T)$$

All $P(T)P(W|T)$ can now be computed

24

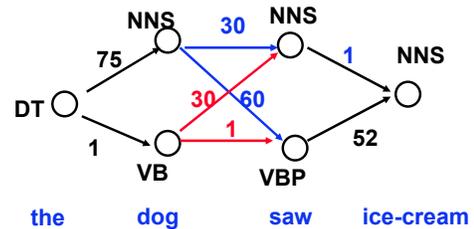
Example



How do we find best path?

25

The counts add scores - we want to find the maximum scoring path



26

How do we find maximum (best) path?

- We use best-first (A*) search, as in AI...
- 1. At each step, k best values (\hat{T}) are chosen. Each of the k values corresponds to one possible tagging combination of the visited words.
- 2. When tagging the next word, recompute probabilities. Go to step 1.
- **Advantage:** fast (do not need to check all possible combinations, but only k potential ones).
- **Disadvantage:** may not return the best solution, but only acceptable results.

27

Accuracy

- Accuracy of this method > **96%**
- **Baseline? 90%**
 - Baseline is performance of stupidest possible method
 - Tag every word with its most frequent tag
 - Tag unknown words as nouns
- Human: **97% +/- 3%**; if discuss together: **100%**

28

Suppose we don't have training data

- Can estimate roughly:
 - start with uniform probabilities,
 - use **Expectation Maximization (EM) algorithm** to re-estimate from counts
 - try labeling with current estimate
 - use this to correct estimate
- Not work well, a small amount of hand-tagged training data improves the accuracy

29

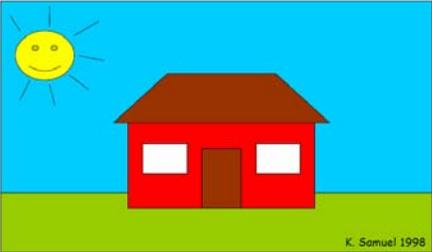
Second approach: transformation-based tagging

Transformation-based Learning (TBL):

- Combines symbolic and stochastic approaches: uses machine learning to refine its tags, via several passes
- Tag using a broadest (most general) rule; then a narrower rule, that changes a smaller number of tags, and so on.

30

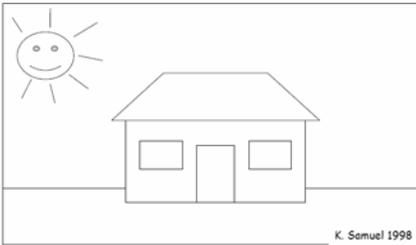
Transformation-based painting



K. Samuel 1998

31

Transformation-based painting



K. Samuel 1998

32

Transformation-based painting



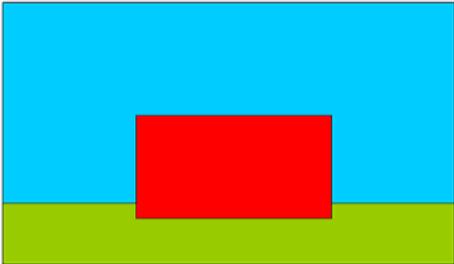
33

Transformation-based painting



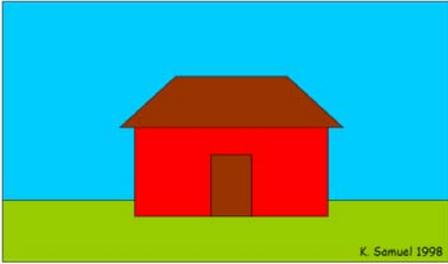
34

Transformation-based painting



35

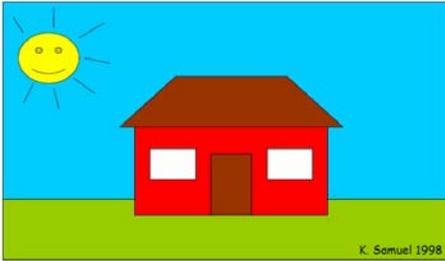
Transformation-based painting



K. Samuel 1998

36

Transformation-based painting



37

How does the TBL system work?

lexicon

data:NN
decided:VB
her:PN
she:PN N
table:NN VB
to:TO

rules

```
pos:NN>VB <- pos:TO@[-1] o
pos:VB>NN <- pos:DT@[-1] o
....
```

input

She decided to table her data
NP VB TO MB PN NN

38

How does the TBL system work?

1. Label every word with its most-likely tag (often 90% right). From Brown corpus:
 $P(\text{NN}|\text{race}) = 0.98$
 $P(\text{VB}|\text{race}) = 0.02$
2. ...expected/VBZ to/(TO race/VB) tomorrow/NN
...the/DT race/NN for/IN outer/JJ space/NN
3. Use transformational (learned) rules:
Change NN to VB when the previous tag is TO
 $\text{pos: 'NN'>'VB'} \leftarrow \text{pos: 'TO' @[-1] o}$

39

Rules for POS tagging

```
pos: 'NN'>'VB' <- pos: 'TO'@[-1] o
pos: 'VBP'>'VB' <- pos: 'MD'@[-1,-2,-3] o
pos: 'NN'>'VB' <- pos: 'MD'@[-1,-2] o
pos: 'VB'>'NN' <- pos: 'DT'@[-1,-2] o
pos: 'VBD'>'VEN' <- pos: 'VBZ'@[-1,-2,-3] o
pos: 'VEN'>'VBD' <- pos: 'PRP'@[-1] o
pos: 'POS'>'VBZ' <- pos: 'PRP'@[-1] o
pos: 'VB'>'VBP' <- pos: 'NNS'@[-1] o
pos: 'IN'>'RB' <- wd:as@[0] & wd:as@[2] o
pos: 'IN'>'WDT' <- pos: 'VB'@[1,2] o
pos: 'VB'>'VBP' <- pos: 'PRP'@[-1] o
pos: 'IN'>'WDT' <- pos: 'VBZ'@[1] o
....
```

40

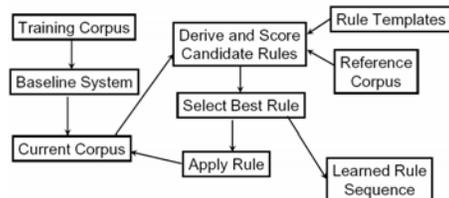
Rules for POS tagging

```

NN VB PREVTAG TO
VB VBP PREVTAG FRP
VBD VEN PREVIOR2TAG VBD
VEN VED PREVTAG FRP
NN VB PREVIOR2TAG MD
VB VBP PREVTAG NNS
VB NN PREVIOR2TAG DT
VEN VED PREVTAG NNP
VBD VEN PREVIOR2OR3TAG VBE
IN DT PREVTAG IN
VBP VB PREVIOR2OR3TAG MD
IN RB WDANDIAFT as as
VED VEN PREVIOR2TAG VB
RE JJ NEXTTAG NN
VBP VB PREVIOR2OR3TAG TO
POS VBZ PREVTAG FRP
NN VBP PREVTAG FRP
DT PDT NEXTTAG DT
...
    
```

41

Learning TB rules in TBL system



Stop when score of best rule falls below threshold.

42

Various Corpora

- Training corpus
w0 w1 w2 w3 w4 w5 w6 w7 w8 w9 w10
- Current corpus (CC 1)
dt vb nn dt vb kn dt vb ab dt vb
- Reference corpus
dt nn vb dt nn kn dt jj kn dt nn

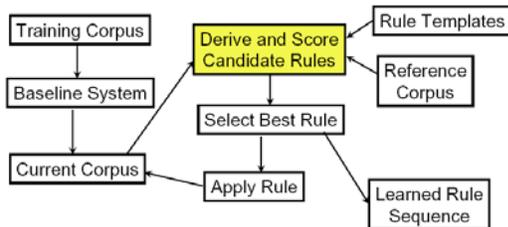
43

Rule Templates

- In TBL, only rules that are instances of *templates* can be learned.
- For example, the rules
tag:'VB'>'NN' ← tag:'DT'@[-1].
tag:'NN'>'VB' ← tag:'DT'@[-1].
are instances of the template
tag:A>B ← tag:C@[-1].
- Alternative syntax using anonymous variables
tag:>_ ← tag:_@[-1].

44

Learning TB rules in TBL system



45

Score, Accuracy and Thresholds

- The *score* of a rule:
$$\text{score}(R) = |\text{pos}(R)| - |\text{neg}(R)|$$
- The *accuracy* of a rule:
$$\text{accuracy}(R) = \frac{|\text{pos}(R)|}{|\text{pos}(R)| + |\text{neg}(R)|}$$
- *Threshold*: the value that a rule must have in order to be considered.
- In *ordinary* TBL, use accuracy threshold < 0.5.

46

Derive and Score Candidate Rule 1

- Template = tag: _>_ ← tag: _@[-1]
- R1 = tag:vb>nn ← tag:dt@[-1]

CC i	dt	vb	nn	dt	vb	kn	dt	vb	ab	dt	vb
CC i+1	dt	nn	nn	dt	nn	kn	dt	nn	ab	dt	nn
Ref. C	dt	nn	vb	dt	nn	kn	dt	jj	kn	dt	nn

- pos(R1) = 3
- neg(R1) = 1
- score(R1) = pos(R1) - neg(R1) = 3-1 = 2

47

Derive and Score Candidate Rule 2

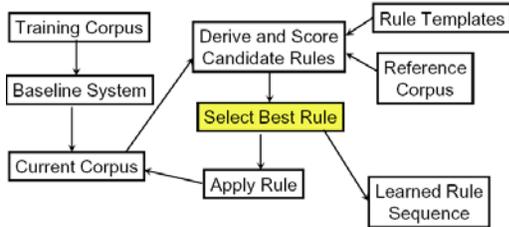
- Template = tag: _>_ ← tag: _@[-1]
- R2 = tag:nn>vb ← tag:vb@[-1]

CC i	dt	vb	nn	dt	vb	kn	dt	vb	ab	dt	vb
CC i+1	dt	vb	vb	dt	vb	kn	dt	vb	ab	dt	vb
Ref. C	dt	nn	vb	dt	nn	kn	dt	nn	kn	dt	nn

- pos(R2) = 1
- neg(R2) = 0
- score(R2) = pos(R2) - neg(R2) = 1-0 = 1

48

Learning TB rules in TBL system



Stop when score of best rule falls below threshold.

49

Select Best Rule

- Current ranking of rule candidates
 $R1 = \text{tag:vb} > \text{nn} \leftarrow \text{tag:dt} @ [-1]$ Score = 2
 $R2 = \text{tag:nn} > \text{vb} \leftarrow \text{tag:vb} @ [-1]$ Score = 1
 ...
- If score threshold ≤ 2 then select R1
- else if score threshold > 2 , terminate.

50

Select Best Rule Optimizations

- **Reduce redundance rules:** only generate candidate rules that have at least one match in the training data.
- **Incremental evaluation:**
 - Keep track of the leading rule candidate.
 - Ignore rules that has #positive matches $<$ score of the leading rule

51

Greedy Best-First Search

Evaluation function

$h(n)$ = estimated cost of the cheapest path from the state represented by the node n to a goal state

52

Advantages of TB Tagging

- Rules can be created/edited manually
- Rules have a declarative, logical semantics
- Simple to implement
- Can be extremely fast (but implementation is more complex)

53

Error analysis: what's hard for taggers

Common errors ($> 4\%$)

- NN (common noun) vs .NNP (proper noun) vs. JJ (adjective): hard to distinguish; important to distinguish especially for information extraction
- RP (particle) vs. RB (adverb) vs. IN (preposition): all can appear in sequences immediate after verb
- VBD vs. VBN vs. JJ: distinguish past tense, past participles, adjective (*raced vs. was raced vs. the out raced horse*)

54

Most powerful unknown word detectors



- 3 inflectional endings (-ed, -s, -ing); 32 derivational endings (-ion, etc.); capitalization; hyphenation
- More generally:
 - Morphological analysis
 - Machine learning approaches