

# Spelling errors

- Assumption: the only spelling errors are a single insertion, deletion, substitution or transposition:

## Example: Spelling errors

insertion: the → ther

deletion: the → th

substitution: the → thw

transposition: the → hte

- The **noisy channel model**: the surface form is an instance of the lexical form which has been passed through a noisy communication channel.
- Spelling error detection has to restore the original form from the noisy instance.

# Spelling errors

## Example: Spelling errors

Error	Correction	Correct	Error	Position	Type
acress	actress	t	ε	2	deletion
acress	cress	ε	a	0	insertion
acress	caress	ca	ac	0	transposition
acress	access	c	r	2	substitution
acress	across	o	e	3	substitution
acress	acres	ε	s	5	insertion
acress	acres	ε	s	4	insertion

# Bayesian inference

- Spelling correction as a *classification* problem: given an *observation* (misspelled word), determine which of a set of *classes* (correctly spelled words) it belongs to.
- Given a vocabulary  $V$  and an observation  $O$ , the (estimated) correct word  $\hat{w}$  is:

$$\hat{w} = \arg \max_{w \in V} P(w|O)$$

- The problem: how to (directly) compute  $P(w|O)$ .

# Bayesian inference

- Bayes rule:

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

- Hence,

$$\begin{aligned}\hat{w} &= \arg \max_{w \in V} P(w|O) \\ &= \arg \max_{w \in V} \frac{P(O|w)P(w)}{P(O)} \\ &= \arg \max_{w \in V} P(O|w)P(w)\end{aligned}$$

because  $P(O)$  is independent of  $w$ , and we are maximizing over all words.

- $P(O|w)$  is the **likelihood**, and  $P(w)$  is the **prior probability**.

## Spelling errors

- Let  $O$  be a misspelled word and  $V$  a set of corrections. Then the most likely correction is:

$$\hat{w} = \arg \max_{w \in V} P(O|w)P(w)$$

- The prior probability of each correction,  $P(w)$ , can be estimated from a corpus by counting how many times  $w$  occurs in the corpus  $\#(w)$  and normalizing by the size of the corpus,  $N$ :

$$P(w) \approx \frac{\#(w)}{N}$$

- Zero counts can cause problems, and hence we **smooth**:

$$P(w) = \frac{\#(w) + 0.5}{N + 0.5V}$$

## Spelling errors: estimating the prior

### Example: Prior probabilities

In a particular corpus of  $N = 44$  million words, the following data were observed:

$w$	$\#(w)$	$P(w)$
actress	1343	0.0000315
cress	0	0.000000014
caress	4	0.0000001
access	2280	0.000058
across	8436	0.00019
acres	2879	0.000065

## Spelling errors: estimating the likelihood

- How to estimate  $P(O|w)$ , the probability of a typo given the correct word?
- The exact probability depends on various factors (who the typist is, etc.)
- Factors which can be estimated include the identity of the letters (e.g., **m** is substituted for **n** because their pronunciation is similar and because they are next to each other on the keyboard) and on context (because they are pronounced similarly, they occur in similar contexts).
- A simplification: using a **confusion matrix** which specifies the number of times one letter was substituted for another in a corpus of errors.

## Spelling errors: estimating the likelihood

### Example: Confusion matrices

$\text{del}(x, y)$ : The number of times the characters  $xy$  were typed as  $x$

$\text{ins}(x, y)$ : The number of times the character  $x$  was typed as  $xy$

$\text{sub}(x, y)$ : The number of times the character  $x$  was typed as  $y$

$\text{trans}(x, y)$ : The number of times the characters  $xy$  were typed as  $yx$ .

$$P(O|w) = \begin{cases} \frac{\text{del}(w_{p-1}, w_p)}{\text{count}(w_{p-1} w_p)} & \text{if deletion} \\ \frac{\text{ins}(w_{p-1}, O_p)}{\text{count}(w_{p-1})} & \text{if insertion} \\ \frac{\text{sub}(O_p, w_p)}{\text{count}(w_p)} & \text{if substitution} \\ \frac{\text{trans}(w_p, w_{p+1})}{\text{count}(w_p w_{p+1})} & \text{if transposition} \end{cases}$$



## Spelling errors: putting everything together

Example: Ranking of candidate corrections for *acress*

<i>w</i>	$\#(w)$	$P(w)$	$P(O w)$	$P(w)P(O w)$	%
<i>actress</i>	1343	0.0000315	0.000117	$3.69^{-9}$	37
<i>cress</i>	0	0.000000014	0.00000144	$2.02^{-14}$	0
<i>caress</i>	4	0.0000001	0.00000164	$1.64^{-13}$	0
<i>access</i>	2280	0.000058	0.000000209	$1.21^{-11}$	0
<i>across</i>	8436	0.00019	0.0000093	$1.77 \times 10^{-9}$	18
<i>acres</i>	2879	0.000065	0.0000321	$2.09 \times 10^{-9}$	21
<i>acres</i>	2879	0.000065	0.0000342	$2.22 \times 10^{-9}$	23

Results: *acres* (normalized percentage of 45%), *actress* (37%).

# Minimum edit distance

- Motivation: The previous (over-simplifying) method assumed that each word had only a single error.
- In general, the problem is that of finding the *distance* between two strings.
- **Minimum edit distance:** the minimum number of operations (insert, delete or substitute) needed to transform one string into another.
- **Levenshtein distance:** each operation has the same cost. Variant: insertions and deletions cost 1, substitutions not allowed.
- Variant: assign a cost to each instance of the operations, e.g., using confusion matrices.

# Minimum edit distance

## Example: Computing minimum edit distance

For two strings,  $s$  and  $t$ ,

$$\text{dist}(i, j) = \begin{cases} 0 & \text{if } i = j = 0 \\ i & \text{if } j = 0 \\ j & \text{if } i = 0 \\ \min \left\{ \begin{array}{l} \text{dist}(i - 1, j) + \text{ins-cost}(t_j), \\ \text{dist}(i - 1, j - 1) + \text{subst-cost}(s_j, t_j), \\ \text{dist}(i, j - 1) + \text{del-cost}(s_j) \end{array} \right\} & \text{otherwise} \end{cases}$$

# Part of speech tagging

- The problem: given a sentence  $O = o_1, \dots, o_n$ , assign to each word  $o_i$  a correct part of speech (POS)  $t_i$
- Resources:
  - **Tagset** a set of POS tags
  - **Lexicon** a list of words with associated possible POS tags
  - **Training data** a corpus where each word is correctly tagged
- Tagsets for English vary, but most have 40–150 tags
- Why is this task important?
- POS tagging for languages with complex morphology...

# POS tagging

## Example: The Penn Treebank Tagset

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &amp;</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>	WP\$	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>( [ , { , &lt;</i>
PP\$	Possessive pronoun	<i>your, one's</i>	)	Right parenthesis	<i>( ], , } , &gt;</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>			

## Part of speech tagging

Example: POS ambiguity in two corpora

English		Hebrew	
#Tags	#word types	#analyses	#tokens
1	35,340	1	17,106
2	3,760	2	7,753
3	264	3	4,619
4	61	4	3,272
5	12	5	1,512
6	2	6	648
7	1	7	388
		8	307
		9	128
		10-12	78

# Markov Model POS tagging

- Assumptions:

**limited horizon:** the tag of a word depends only on the tag of the previous word

**time invariant:** this dependency does not change over time

For example, if a pronoun has a probability  $p$  to occur after an auxiliary verb in the beginning of a sentence, then this probability does not change in the rest of the sentence.

- Plausibility?

- Notation: subscripts refer to positions in the sentence and in the corpus; superscripts refer to word types in the lexicon and tag types in the tagset.

## Markov Model POS tagging

- The states of the Markov model are tags; each time the computation leaves a state, a word is emitted
- The maximum likelihood estimates of some tag  $t^k$  following a tag  $t^j$  are estimated from the tags' relative frequencies:

$$P(t^k|t^j) = \frac{\#(t^j t^k)}{\#(t^j)}$$

This constitutes the values of the transition probabilities  $a_{ij}$

- The probability of a word being emitted by a particular state (tag) via maximum likelihood estimation:

$$P(w^l|t^j) = \frac{\#(w^l : t^j)}{\#(t^j)}$$

This constitutes the values of the emission probabilities  $b_{ijk}$ .



## Markov Model POS tagging

- The best tagging  $t_{1..n}$  for a sentence  $w_{1..n}$  is:

$$\begin{aligned}\arg \max_{t_{1..n}} P(t_{1..n}|w_{1..n}) &= \arg \max_{t_{1..n}} \frac{P(w_{1..n}|t_{1..n})P(t_{1..n})}{P(w_{1..n})} \\ &= \arg \max_{t_{1..n}} P(w_{1..n}|t_{1..n})P(t_{1..n})\end{aligned}$$

- Assuming (wrongly!) that words are independent of each other, and that a word's identity depends only on its tag,

$$\begin{aligned}P(w_{1..n}|t_{1..n}) &= \prod_{i=1}^n P(w_i|t_{1..n}) \\ &= \prod_{i=1}^n P(w_i|t_i)\end{aligned}$$

- By partitioning and the assumption of limited horizon,

$$\begin{aligned}P(t_{1..n}) &= P(t_n|t_{1..n-1})P(t_{n-1}|t_{1..n-2})\cdots P(t_2|t_1) \\ &= P(t_n|t_{n-1})P(t_{n-1}|t_{n-2})\cdots P(t_2|t_1)\end{aligned}$$

# Markov Model POS tagging

- The best tagging  $t_{1..n}$  for a sentence  $w_{1..n}$  is:

$$\arg \max_{t_{1..n}} P(t_{1..n}|w_{1..n}) = \prod_{i=1}^n P(w_i|t_i)P(t_i|t_{i-1})$$

- A direct evaluation would require an exponential number of multiplications
- Use the Viterbi algorithm for classification:
  - $\delta[t, i]$  is the probability of being at state (tag)  $i$  at time (word)  $t$
  - $\Psi[t + 1, i]$  is the most likely state (tag) at time (word)  $t$  given that at time (word)  $t$  we are in state (tag)  $i$

# Markov Model POS tagging

## Example: Viterbi classification

Initialization:

$$\delta[., 1] = 1.0; \delta[t, 1] = 0.0 \text{ for } t \neq .;$$

Induction:

for  $i = 1$  to  $n$

  for all tags  $t^j$

$$\delta[t^j, i + 1] = \max_{1 \leq k \leq T} (\delta[t^k, i] P(w_{i+1} | t^j) P(t^j | t^k));$$

$$\Psi[t^j, i + 1] = \arg \max_{1 \leq k \leq T} (\delta[t^k, i] P(w_{i+1} | t^j) P(t^j | t^k));$$

Termination and path read-out:

$$t_{n+1} = \arg \max_{1 \leq j \leq T} \delta[j, n + 1];$$

for  $j = n$  downto 1

$$t_j = \Psi[t_{j+1}, j + 1];$$

$$P(t_1, \dots, t_n) = \max_{1 \leq j \leq T} \delta[t^j, n + 1];$$

## Shallow parsing

- **Shallow parsing** consists of identifying the main components of sentences and their heads and determining syntactic relationships among them.
- **Problem:** Given an input string  $O = \langle o_1, \dots, o_n \rangle$ , a **phrase** is a consecutive substring  $\langle o_i, \dots, o_j \rangle$ . The goal is, given a sentence, to identify all the phrases in the string.
- A secondary goal is to tag the phrases as Noun Phrase, Verb Phrase etc.
- An additional goal is to identify relations between phrases, such as subject-verb, verb-object etc.
- **Question:** How can this problem be cast as a classification problem?

## Text chunking

- Text chunking involves dividing sentences into non-overlapping segments on the basis of fairly simple superficial analysis.
- This is a useful and relatively tractable precursor to full parsing, since it provides a foundation for further levels of analysis, while still allowing complex attachment decisions to be postponed to a later phase.

## Deriving chunks from treebank parses

- Annotation of training data can be done automatically based on the parsed data of the Penn Tree Bank
- Two different chunk structure tagsets: one bracketing non-recursive “base NPs”, and one which partitions sentences into non-overlapping N-type and V-type chunks

## “Base NP” chunk structure

- The goal of the “base NP” chunks is to identify essentially the initial portions of non-recursive noun phrases up to the head, including determiners but not including postmodifying prepositional phrases or clauses.
- These chunks are extracted from the Treebank parses, basically by selecting NPs that contain no nested NPs.
- The handling of conjunction follows that of the Treebank annotators as to whether to show separate baseNPs or a single baseNP spanning the conjunction.
- Possessives are treated as a special case, viewing the possessive marker as the first word of a new baseNP, thus flattening the recursive structure in a useful way.

## “Base NP” chunk structure

Example: “Base NP” chunk structure

[*N* The government *N*] has [*N* other agencies and instruments *N*] for pursuing [*N* these other objectives *N*] .

Even [*N* Mao Tse-tung *N*] [*N* 's China *N*] began in [*N* 1949 *N*] with [*N* a partnership *N*] between [*N* the communists *N*] and [*N* a number *N*] of [*N* smaller , non-communist parties *N*] .



## Partitioning chunks

- In the partitioning chunk experiments, the prepositions in prepositional phrases are included with the object NP up to the head in a single N-type chunk.
- The handling of conjunction again follows the Treebank parse.
- The portions of the text not involved in N-type chunks are grouped as chunks termed V-type, though these “V” chunks include many elements that are not verbal, including adjective phrases.
- Again, the possessive marker is viewed as initiating a new N-type chunk.

# Partitioning chunks

## Example: Partitioning chunks

[N Some bankers N] [V are reporting V] [N more inquiries than usual  
N] [N about CDs N] [N since Friday N] .

[N Indexing N] [N for the most part N] [V has involved simply buying  
V] [V and then holding V] [N stocks N] [N in the correct mix N] [V  
to mirror V] [N a stock market barometer N] .

## Encoding chunking as a tagging problem

- Each word carries both a part-of-speech tag and also a “chunk tag” from which the chunk structure can be derived.
- In the baseNP experiments, the chunk tag set is  $\{I, O, B\}$ , where words marked *I* are *inside* some baseNP, those marked *O* are *outside*, and the *B* tag is used to mark the leftmost item of a baseNP which immediately follows another baseNP.
- In the partitioning experiments, the chunk tag set is  $\{BN, N, BV, V, P\}$ , where *BN* marks the first word and *N* the succeeding words in an N-type group while *BV* and *V* play the same role for V-type groups.

## Encoding chunking as a tagging problem

- Encoding chunk structure with tags attached to words (rather than inserting bracket markers between words) limits the dependence between different elements of the encoded representation.
- While brackets must be correctly paired in order to derive a chunk structure, it is easy to define a mapping that can produce a valid chunk structure from any sequence of chunk tags; the few hard cases that arise can be handled locally.
- For example, in the baseNP tag set, whenever a *B* tag immediately follows an *O*, it must be treated as an *I*.
- In the partitioning chunk tag set, wherever a *V* tag immediately follows an *N* tag without any intervening *BV*, it must be treated as a *BV*.

## Transformation-based learning for Chunking

- Transformational learning begins with some initial “baseline” prediction, which here means a baseline assignment of chunk tags to words.
- Reasonable suggestions for baseline heuristics after a text has been tagged for part-of-speech might include assigning to each word the chunk tag that it carried most frequently in the training set, or assigning each part-of-speech tag the chunk tag that was most frequently associated with that part-of-speech tag in the training.
- Testing both approaches, the baseline heuristic using part-of-speech tags turned out to do better.

## Rule templates

- Rules can refer to words and to POS tags. Up to three words to the left and right of the target word, and up to two POS tags to the left and right of the target can be addressed.
- A set of 100 rule templates, obtained by the cross product of 20 word-patterns and 5 tag-patterns, was used.
- Then, a variant of Brill's TBL algorithm was implemented.

# Results

- BaseNP chunks:

Training	Recall	Precision
Baseline	81.9%	78.2%
50K	90.4%	89.8%
100K	91.8%	91.3%
200K	92.3%	91.8%

- Partitioning chunks:

Training	Recall	Precision
Baseline	60.0%	47.8%
50K	86.6%	85.8%
100K	88.2%	87.4%
200K	88.5%	87.7%

## Memory-based shallow parsing

- Shallow parsing consists of discovering the main constituents of sentences (NPs, VPs, PPs) and their heads, and determining syntactic relationships (like subjects, objects or adjuncts) between verbs and heads of other constituents.
- This is an important component of text analysis systems in applications such as information extraction and summary generation.



## Memory-based learning: reminder

- A memory-based learning algorithm constructs a classifier for a task by storing a set of examples.
- Each example associates a feature vector (the problem description) with one of a finite number of classes (the solution).
- Given a new feature vector, the classifier extrapolates its class from those of the most similar feature vectors in memory.
- The metric defining similarity can be automatically adapted to the task at hand.

# Organization

- Syntactic analysis is carved up into a number of classification tasks.
- These can be segmentation tasks (e.g., deciding whether a word or tag is the beginning or the end of an NP) or disambiguation tasks (e.g., deciding whether a chunk is the subject, object or neither).
- Output of one module (e.g., POS tagging or chunking) is used as input by other modules (e.g., syntactic relation assignment).

# Algorithms and implementation

- All the experiments use TiMBL.
- Two variants of MBL are used:
  - IB1-IG: The distance between a test item and a memory item is the number of features on which they disagree. The algorithm uses information gain to weigh the cost of mismatches. Classification speed is linear in the number of training instances times the number of features.
  - IGTREE: A decision tree is created with features as tests, ordered according to information gain of features. Classification speed is linear in the number of features times the average branching factor of the tree, which is bound by the average number of values per feature.

# Experiments

- Two series of experiments:
  - Memory-based NP and VP chunking
  - Subject/object detection using the chunker

# Chunking as a tagging task

- Each word is assigned a tag which indicates whether it is inside or outside a chunk:
  - I\_NP inside a baseNP
  - O outside both a baseNP and a baseVP
  - B\_NP inside a baseNP, but the preceding word is in another baseNP
  - I\_VP inside a baseVP
  - B\_VP inside a baseVP, but the preceding word is in another baseVP
- Since baseNPs and baseVPs are non-overlapping and non-recursive, these five tags suffice to unambiguously chunk a sentence.

# Tagging example

Example:

[NP Pierre<sub>I\_NP</sub> Vinken<sub>I\_NP</sub> NP] ,O [NP 61<sub>I\_NP</sub> years<sub>I\_NP</sub> NP] old<sub>O</sub> ,O  
[VP will<sub>I\_VP</sub> join<sub>I\_VP</sub> VP] [NP the<sub>I\_NP</sub> board<sub>I\_NP</sub> NP] as<sub>O</sub> [NP a<sub>I\_NP</sub>  
nonexecutive<sub>I\_NP</sub> director<sub>I\_NP</sub> NP] [NP Nov.<sub>B\_NP</sub> 29<sub>I\_NP</sub> NP] .

## Chunking as a tagging task: experiments

- The features for the experiments are the word and the POS tag (as provided by the Penn Tree Bank) of two words to the left, the target word and one word to the right.
- The baseline is computed with IB1-IG, using as features only the focus word/POS.

# Results

- BaseNP chunks:

Method	Recall	Precision
Baseline words	79.7%	76.2%
Baseline POS	82.4%	79.5%
IGTree	93.1%	91.8%
IB1-IG	94.0%	93.7%

- BaseVP chunks:

Method	Recall	Precision
Baseline words	73.4%	67.5%
Baseline POS	87.7%	74.7%
IGTree	94.2%	93.0%
IB1-IG	95.5%	94.0%



## Subject/object detection

- Finding the subject or object of a verb is defined as a mapping from pairs of words (the verb and the head of the constituent), and a representation of their context, to a class (subject, object or neither).
- A verb can have multiple subjects (in the case of NP coordination) and a word can be the subject of more than one verb (VP coordination).
- The input is POS tagged and chunked.

### Example:

[*NP* My/*PRP* sisters/*NNS* *NP*] [*VP* have/*VBP* not/*RB* seen/*VCN*  
*VP*] [*NP* the/*DT* old/*JJ* man/*NN* *NP*] lately/*RB* ./.

- All chunks are reduced to their heads, defined as the rightmost word of a baseNP or baseVP.

## Subject/object detection: features

- The distance, *in chunks*, between the verb and the head
- The number of other baseVPs between the verb and the head
- The number of commas between the verb and the head
- The verb and its POS tag
- The head and its POS tag
- The two left context and one right context words/chunks of the head, represented by the word and its POS tag

## Subject/object detection: results

Finding unrestricted subjects and objects is hard:

Method	Together	Subjects	Objects
Heuristic baseline	66.2	65.2	67.7
IGTree	76.2	75.8	76.8
IB1-IG	75.6	76.5	74.0
Unanimous	77.8	77.1	79.0

## The use of classifiers in sequential inference

- Combination of the outcome of several classifiers in a way that provides a coherent inference that satisfies some constraints.
- Two general approaches to identifying phrase structure: *projection-based Markov models* and *constraint satisfaction with classifiers*.

## Identifying phrase structure

- Classifiers recognize in the input string local signals which are indicative of the existence of phrases
- Classifiers can indicate that an input symbol is **i**nside or **o**utside a string or that a symbol **o**pens or **c**loses a string.
- The open/close approach has been found more robust and is pursued here.
- The classifiers' outcomes can be combined to determine the phrase, but this combination must satisfy certain constraints for the result to be legitimate.
- Several types of constraints, such as length, order and others, can be formalized and incorporated into the two approaches studied here.

## Identifying phrase structure

- Two complex phrase identification tasks are defined: base NPs and Subject-Verb patterns.

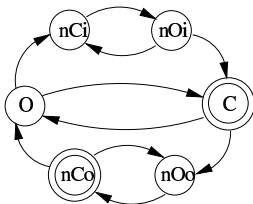
### Example:

[ The theory presented claims ] that [ the algorithm runs ] and performs ...

- Two classifiers are learned for each task, predicting whether word  $t$  **o**pens or **c**loses a phrase.

## Identifying phrase structure

- Each classifier may output two values: open/ $\neg$ open and close/ $\neg$ close.
- However, for technical reasons, *three* values are output by each classifier, where the 'not' value is divided according to whether or not the word is inside a phrase.
- Consequently, the values are: **O**, **nOi**, **nOo**, **C**, **nCi**, **nCo**.
- The order of these values is constrained according to the following diagram:



# Definitions

- The input string is  $O = \langle o_1, o_2, \dots, o_n \rangle$
- A phrase  $\pi^{i,j}(O)$  is a substring  $\langle o_1, o_{i+1}, \dots, o_j \rangle$  of  $O$
- $\pi^*(O)$  is the set of all possible phrases of  $O$
- $\pi^{i,j}(O)$  and  $\pi^{k,l}(O)$  **overlap**, denoted  $\pi^{i,j}(O) \rightleftharpoons \pi^{k,l}(O)$ , iff  $j \geq k$  and  $l \geq i$
- Given a string  $O$  and a set  $Y$  of classes of phrases, a solution to the phrase identification problem is a set  $\{(\pi, y) \mid \pi \in \pi^*(O) \text{ and } y \in Y\}$  such that for all  $(\pi_i, y), (\pi_j, y)$ , if  $i \neq j$  then  $\pi_i \neq \pi_j$
- We assume that  $|Y| = 1$ .



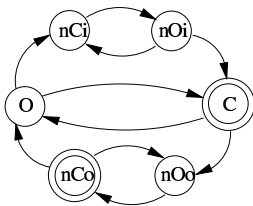
# Hidden Markov Model combinator

Reminder: an HMM is a probabilistic finite state automaton consisting of

- A finite set  $S$  of states
- A set  $O$  of observations
- An initial state distribution  $P_1(s)$
- A state-transition distribution  $P(s|s')$  for  $s, s' \in S$ , and
- An observation distribution  $P(o|s)$  for  $o \in O, s \in S$ .
- Constraints can be incorporated into the HMM by constraining the state transition probability distribution. For example, set  $P(s|s') = 0$  for cases where the transition from  $s'$  to  $s$  is not allowed.

## Hidden Markov Model combinator

- We assume that we have local signals which indicate the state. That is, classifiers are given with states as their outcome.
- Formally, we assume that  $P_t(s|o_y)$  is given, where  $t$  is a time step in the sequence.
- 



- Constraints on state transitions do not have to be stated explicitly; they can be recovered from training data.

## Hidden Markov Model combinator

- Instead of estimating the observation probability  $P(o|s)$  directly from training data, it is computed from the classifiers' output:

$$P_t(o_t|s) = \frac{P_t(s|o_t) \times P_t(o_t)}{P_t(s)}$$

- 

$$P_t(s) = \sum_{s' \in S} P(s|s') \times P_{t-1}(s')$$

where  $P_1(s)$  and  $P(s|s')$  are the standard HMM distributions

- $P_t(o_t)$  can be treated as constant since the observation sequence is fixed for all compared sequences.
- The Viterbi algorithm can be used to find the most likely state sequence for a given observation.

## Projection based Markov Model combinator

- In standard HMMs, observations are allowed to depend only on the current state; no long-term dependencies can be modeled.
- Similarly, constraint structure is restricted by having a stationary probability distribution of a state given the previous state.
- In PMM, these limitations are relaxed by allowing the distribution of a state to depend, in addition to the previous state, on the observation.
- Formally, the independence assumption is:

$$P(s_t | s_{t-1}, \dots, s_1, o_{t-1}, \dots, o_1) = P(s_t | s_{t-1}, o_t)$$

## Projection based Markov Model combinator

- Given an observation sequence  $O$ , the most likely state sequence  $S$  given  $O$  is obtained by maximizing

$$\begin{aligned} P(S|O) &= \prod_{t=2}^n [P(s_t|s_1, \dots, s_{t-1}, o)] P_1(s_1|o) \\ &= \prod_{t=2}^n [P(s_t|s_{t-1}, o_t)] P_1(s_1|o_1) \end{aligned}$$

- In this model the classifiers decisions are incorporated in the terms  $P(s|s', o)$  and  $P_1(s|o)$ . The classifiers take into account not only the current input symbol but also the previous state. The hope is that these new classifiers perform better because they are given more information.

## Constraint satisfaction based combinator

- A boolean constraint satisfaction problem consists of a set of  $n$  variables  $V = \{v_1, \dots, v_n\}$ , each ranging over values in a domain  $D_i$  (here, 0/1).
- A constraint is a relation over a subset of the variables, defining a set of “global” possible assignments to the referred variables.
- A solution to a CSP is an assignment that satisfies all the constraints.
- The CSP formalism is extended to deal with probabilistic variables; the solution now has to minimize some cost function. Thus, each variable is associated with a cost function  $c_i : D_i \rightarrow \mathfrak{R}$ .

## Constraint satisfaction with classifiers

- Given an input  $O = \langle o_1, \dots, o_l \rangle$ , let  $V = \{v_{i,j} \mid 1 \leq i \leq j \leq l\}$ . Each variable  $v_{i,j}$  corresponds to a potential phrase  $\pi^{i,j}(O)$ .
- Associate with each variable  $v_{i,j}$  a cost function  $c_{i,j}$ .
- Constraints can now be expressed as boolean formulae. For example, the constraint that requires that no two phrases overlap is expressed as:

$$\bigwedge_{\pi^{a,b} \cap \pi^{c,d}} (\neg v_{a,b} \vee \neg v_{c,d})$$

- The solution is an assignment of 0/1 to variables which satisfies the constraints and, in addition, minimizes the overall cost

$$\sum_{i=1}^n c_i(v_i)$$

## Constraint satisfaction with classifiers

- In general, the corresponding optimization problem is NP-hard.
- In the special case where costs are in  $[0, 1]$ , a solution which is at most twice the optimal can be found efficiently.
- For the specific case of non-overlapping phrase identification, the problem can be solved efficiently using a graph representation of the constraints (the problem reduces to finding a shortest path in a weighted graph).
- What is left is determining the cost function.



## Constraint satisfaction with classifiers: cost function

- It can be shown that in order to maximize the number of correct phrases, each phrase has to be assigned a cost that is minus the probability of the phrase being correct:

$$c_{i,j}(v_{i,j}) = \begin{cases} -p_{i,j} & \text{if } v_{i,j} = 1 \\ 0 & \text{otherwise} \end{cases}$$

where  $p_{i,j}$  is the probability that the phrase  $\pi^{i,j}$  is correct.

## Constraint satisfaction with classifiers: cost function

- Assuming independence between symbols in a phrase, and assuming that the important part of a phrase are only its beginning and end words,

$$p_{i,j} = P_i^O(O) \times P_j^C(C)$$

where  $P_i^O(O)$  is the probability that the first symbol  $o_i$  in the phrase is actually the beginning of a phrase, and  $P_j^C(C)$  is the probability that the last symbol  $o_j$  of the phrase is actually the end of a phrase.

- These two probabilities are supplied by the classifiers.

# Results

Method	NP		SV	
	POS only	POS + word	POS only	POS + word
HMM	90.64	92.89	64.15	77.54
PMM	90.61	92.98	74.98	86.07
CSC	90.87	92.88	85.36	90.09