

6.891: Lecture 13 (October 22nd, 2003)

Machine Translation Part IV

Announcements

- Philipp Koehn talk tomorrow:

Advances in Statistical Machine Translation: Phrases, Noun Phrases and Beyond

Speaker: Philipp Koehn

Speaker Affiliation: Information Sciences Institute / Univ. of Southern California

Date: 10-23-2003 Time: 2:30 PM - 3:30 PM Refreshments: 2:15 PM Location: NE43-518 / 200 Technology Square (Room 518)

I will review the state of the art in statistical machine translation (SMT), present my dissertation work, and sketch out the research challenges of syntactically structured statistical machine translation.

The best methods currently in SMT build on the translation of phrases (any sequence of words) instead of single words. Phrase translation pairs are automatically learned from parallel corpora. While SMT systems generate translation output that often conveys a lot of the meaning of the original text, it is frequently ungrammatical and incoherent.

The research challenge at this point is to introduce syntactic knowledge to the state of the art in order to improve translation quality. My approach breaks up the translation process along linguistic lines. I will present my thesis work on noun phrase translation and ideas about clause structure.

- I'll post an announcement about projects before the next lecture: please arrange a meeting with me to discuss possible projects

Overview

- Syntax Based Model 1: (Yamada and Knight 2001)
- Syntax Based Model 2: (Wu 1995)
- A Phrase-Based Model: (Koehn, Och and Marcu 2003)

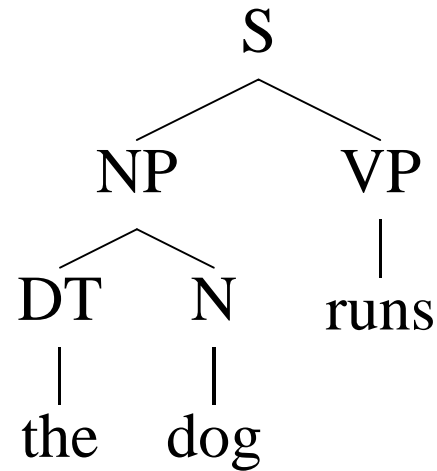
Methods that go beyond word-word alignments

(Yamada and Knight 2001)

- Task: English to Japanese translation
- IBM Models may be poor for languages with very different word orders?
- Task is Japanese \rightarrow English translation, and we have an English parser
- Notation: as before we'll use f as the source language (was French, now Japanese), and e as the target language
- Notation: we'll use \mathcal{E} to refer to an English **tree**

An Example (\mathcal{E} , f) Pair

\mathcal{E} :



f : arun athe adog anow

**Preprocessing of the training set:
Parse all the English strings**

Problems that Need to be Solved

- How to model $P(\mathbf{f} \mid \mathcal{E})$?
i.e., how is a French **string** generated from an English **tree**?

- How do we train the parameters of the model?

- How do we decode with the model, i.e., find

$$\operatorname{argmax}_{\mathbf{e}} P(\mathbf{f} \mid \mathcal{E}) P(\mathbf{e})$$

where \mathbf{e}, \mathcal{E} is a sentence/tree pair in English?

How to model $P(f | e)$?:
Three Operations that Modify Trees

- **Reordering** operations
- **Insertion** of French words
- **Translation** of English words

Reordering Operations

- For each rule with n children, there are $n!$ possible reorderings
- For example, $S \rightarrow \text{ADVP NP VP}$ can be reordered in 6 possible ways

S	→	ADVP	NP	VP
S	→	ADVP	VP	NP
S	→	NP	ADVP	VP
S	→	NP	VP	ADVP
S	→	VP	NP	ADVP
S	→	VP	ADVP	NP

Reordering Operations

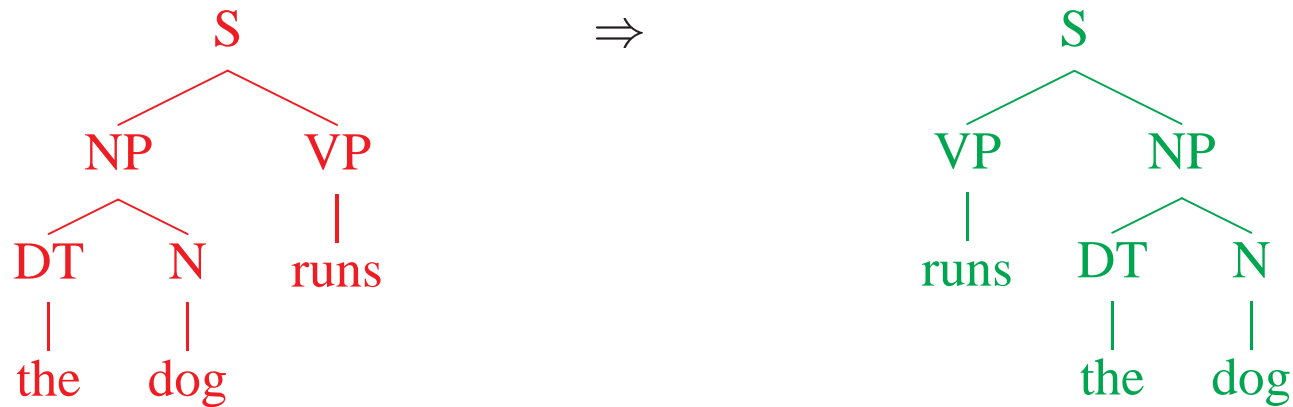
- Introduce $\rho(r' | r)$ as probability of r being reordered as r'
- For example,

$$\rho(S \rightarrow VP \ ADVP \ NP \mid S \rightarrow ADVP \ NP \ VP)$$

- We now have a table of these probabilities for each rule:

r'					$\rho(r' \mid S \rightarrow ADVP \ NP \ VP)$
S	→	ADVP	NP	VP	0.5
S	→	ADVP	VP	NP	0.1
S	→	NP	ADVP	VP	0.3
S	→	NP	VP	ADVP	0.03
S	→	VP	NP	ADVP	0.04
S	→	VP	ADVP	NP	0.03

An Example of Reordering Operations



Has probability:

$$\begin{aligned} & \rho(S \rightarrow VP \ NP \mid S \rightarrow NP \ VP) \times \\ & \rho(NP \rightarrow DT \ N \mid NP \rightarrow DT \ N) \\ & \rho(DT \rightarrow the \mid DT \rightarrow the) \\ & \rho(N \rightarrow dog \mid N \rightarrow dog) \\ & \rho(VP \rightarrow runs \mid VP \rightarrow runs) \end{aligned}$$

Note: Unary rules can only “reorder” in one way, with probability 1
e.g., $\rho(VP \rightarrow runs \mid VP \rightarrow runs) = 1$

Insertion Operations

- At any node in the tree, we can either:
 - Generate no “inserted” foreign words
e.g., has probability

$$\mathbf{I}_1(\text{none} \mid NP, S)$$

here NP is the node in the tree, S is its parent

- Generate an inserted foreign word to the left of the node
e.g., has probability

$$\mathbf{I}_1(\text{left} \mid NP, S)\mathbf{I}_2(\text{anow})$$

here NP is the node in the tree, S is its parent, and anow is inserted to the left of the node

- Generate an inserted foreign word to the right of the node

$$\mathbf{I}_1(\textit{right} \mid NP, S)\mathbf{I}_2(\textit{anow})$$

here NP is the node in the tree, S is its parent, and \textit{anow} is inserted to the right of the node

An Example of Insertion Operations



Has probability:

$\mathbf{I}_1(\text{right} \mid NP, S) \times \mathbf{I}_2(\text{anow}) \times$

$\mathbf{I}_1(\text{none} \mid S, TOP) \times$

$\mathbf{I}_1(\text{none} \mid VP, S) \times$

$\mathbf{I}_1(\text{none} \mid \text{runs}, VP) \times$

$\mathbf{I}_1(\text{none} \mid DT, NP) \times$

$\mathbf{I}_1(\text{none} \mid N, NP) \times$

$\mathbf{I}_1(\text{none} \mid \text{the}, DT) \times$

$\mathbf{I}_1(\text{none} \mid \text{dog}, N)$

Translation Operations

For each English word, translate it to French word f with probability $\mathbf{T}(f \mid e)$ (note that f can be *NULL*)



Has probability:

$$\mathbf{T}(aruns \mid runs) \times \mathbf{T}(athe \mid the) \times \mathbf{T}(adog \mid dog)$$

Summary: Three Operations that Modify Trees

- The three operations:
 - **Reordering** operations with parameters ρ
 - **Insertion** of French words with parameters $\mathbf{I}_1, \mathbf{I}_2$
 - **Translation** of English words with parameters \mathbf{T}
- In this case, the **alignment** \mathbf{a} is the sequence of reordering, insertion and translation operations used to build \mathbf{f}
- We have a model of $P(\mathbf{f}, \mathbf{a} \mid \mathcal{E})$
- Note that each $(\mathcal{E}, \mathbf{f})$ pair may have many possible alignments

- Two questions:

1. How do we train the ρ , \mathbf{I}_1 , \mathbf{I}_2 , \mathbf{T} parameters?
2. How do we find

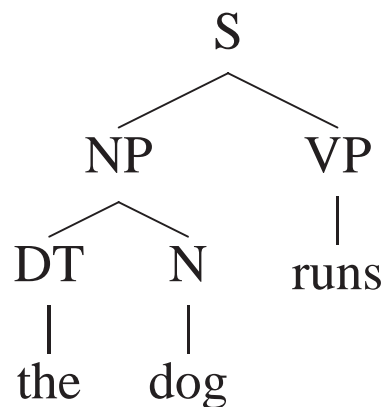
$$\operatorname{argmax}_{\mathcal{E}, \mathbf{e}, \mathbf{a}} P(\mathbf{f}, \mathbf{a} \mid \mathcal{E}) P(\mathbf{e})$$

where $(\mathcal{E}, \mathbf{e}, \mathbf{a})$ is an English tree, sentence, alignment triple?

The translation problem:

Input: arun athe adog anow

Output:



A Slightly Simpler Translation Problem

- For now, instead of trying to find

$$\operatorname{argmax}_{\mathcal{E}, \mathbf{e}, \mathbf{a}} P(\mathbf{f}, \mathbf{a} \mid \mathcal{E}) P(\mathbf{e})$$

we'll consider a method that finds

$$\operatorname{argmax}_{\mathcal{E}, \mathbf{e}, \mathbf{a}} P(\mathbf{f}, \mathbf{a} \mid \mathcal{E})$$

(no language model)

- This can be done by transforming our model into a probabilistic context-free grammar, then parsing the French sentence using dynamic programming!!!

Constructing a PCFG

- For each English/French word pair (e, f) , construct rules

$e \rightarrow f$

with probabilities $\mathbf{T}(f \mid e)$

- For example, $\text{dog} \rightarrow \text{adog}$ with probability $\mathbf{T}(\text{adog} \mid \text{dog})$

- Also construct rules

$e \rightarrow \epsilon$

with probabilities $\mathbf{T}(NULL \mid e)$ (where ϵ is the empty string)

Constructing a PCFG

- For every pair of non-terminals construct rules such as

NP-S \rightarrow NP with probability $\mathbf{I}_1(\text{none} \mid NP, S)$

NP-S \rightarrow INS NP with probability $\mathbf{I}_1(\text{left} \mid NP, S)$

NP-S \rightarrow NP INS with probability $\mathbf{I}_1(\text{right} \mid NP, S)$

- Also, for every French word f that can be inserted, construct rules such as

INS \rightarrow f with probability $\mathbf{I}_2(f)$

e.g.,

INS \rightarrow anow with probability $\mathbf{I}_2(\text{anow})$

Constructing a PCFG

- For every rule in English r , for every reordering of r , construct following rules

(example with $r = S \rightarrow \text{ADVP NP VP}$,
 $r' = S \rightarrow \text{VP ADVP NP}$)

$S \rightarrow S(\text{ADVP}, \text{NP}, \text{VP})$ with probability 1

$S(\text{ADVP}, \text{NP}, \text{VP}) \rightarrow \text{VP-S ADVP-S NP-S}$

with probability $\rho(S \rightarrow \text{VP ADVP NP} \mid S \rightarrow \text{ADVP NP VP})$

Constructing a PCFG

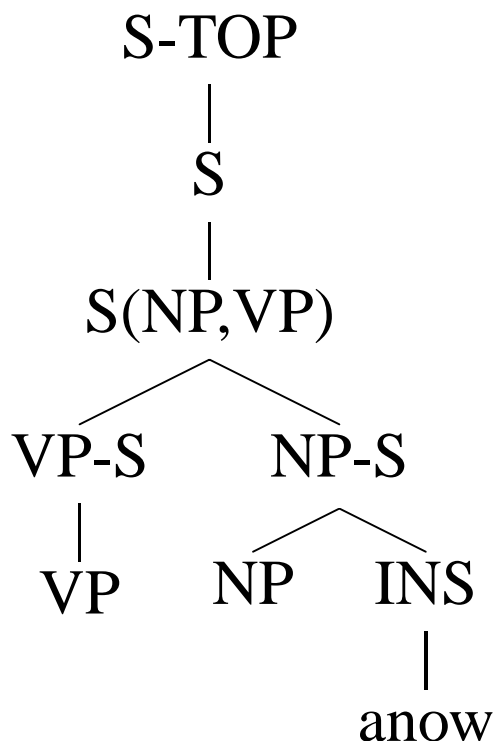
- Finally, for every non-terminal X , construct a start symbol

$X\text{-TOP}$

for example,

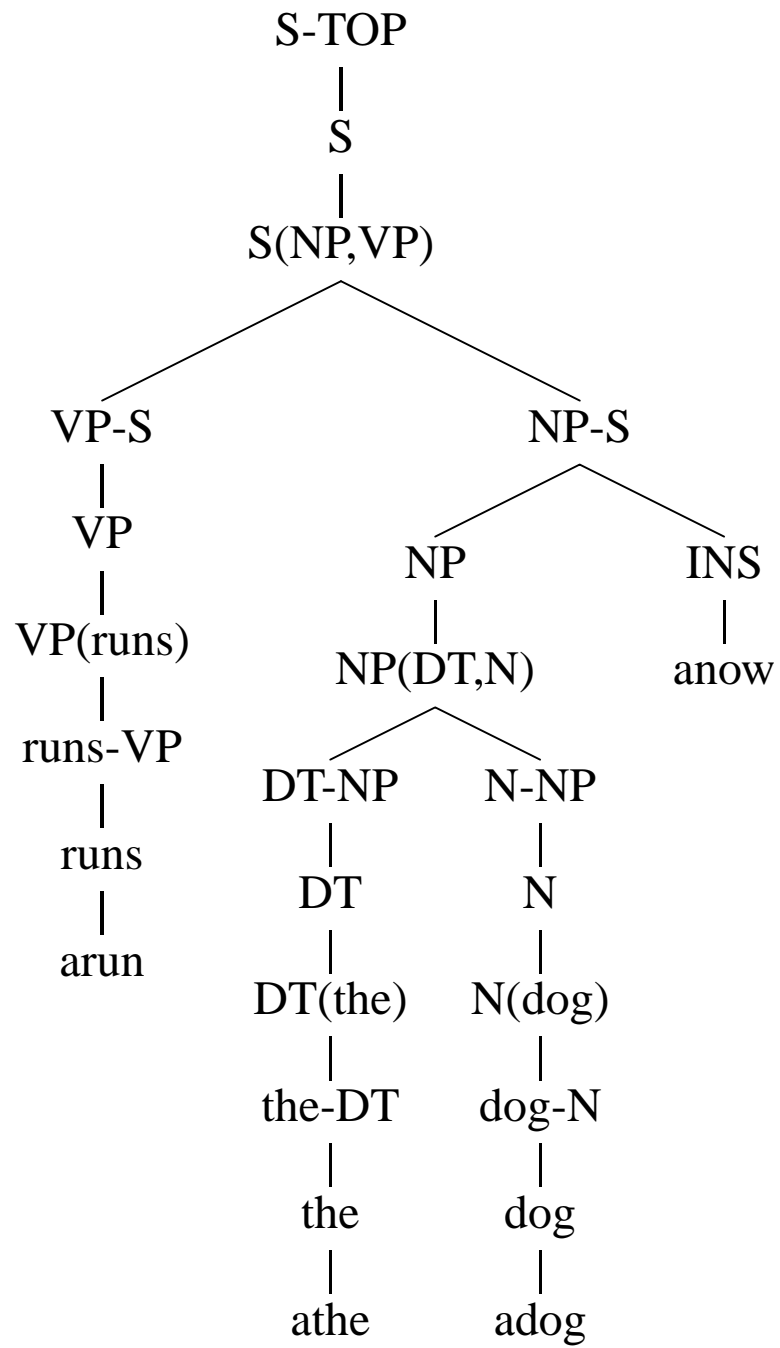
$S\text{-TOP}$

An example:



This subtree has probability:

$$\mathbf{I}_1(\textit{none} \mid S, TOP) \times \rho(S \rightarrow VP \ NP \mid S \rightarrow NP \ VP) \times \\ \mathbf{I}_1(\textit{none} \mid VP, S) \times \mathbf{I}_1(\textit{right} \mid NP, S) \times \mathbf{I}_2(\textit{anow})$$



Other Points

- Once we've constructed the PCFG, finding the most likely parse for a French string → finding the most likely English parse tree, English string, and alignment
- The model can be trained using EM:
dynamic programming approach is possible
- Can parse a French sentence to produce a **forest**:
a compact representation of all possible English translations
- A trigram language model can be used to pick the highest scoring string from the forest (although I'm not sure about the computational complexity of this...)
- (Yamada and Knight 2002) describe newer models

Overview

- Syntax Based Model 1: (Yamada and Knight 2001)
- Syntax Based Model 2: (Wu 1995)
- A Phrase-Based Model: (Koehn, Och and Marcu 2003)

Methods that go beyond word-word alignments

(Wu 1995)

- Standard probabilistic context-free grammars:
probabilities over rewrite rules define probabilities over trees,
strings, in one language
- **Transduction grammars:**
Simultaneously generate strings in two languages

A Probabilistic Context-Free Grammar

S	⇒	NP	VP	1.0
VP	⇒	Vi		0.4
VP	⇒	Vt	NP	0.4
VP	⇒	VP	PP	0.2
NP	⇒	DT	NN	0.3
NP	⇒	NP	PP	0.7
PP	⇒	P	NP	1.0

Vi	⇒	sleeps	1.0
Vt	⇒	saw	1.0
NN	⇒	man	0.7
NN	⇒	woman	0.2
NN	⇒	telescope	0.1
DT	⇒	the	1.0
IN	⇒	with	0.5
IN	⇒	in	0.5

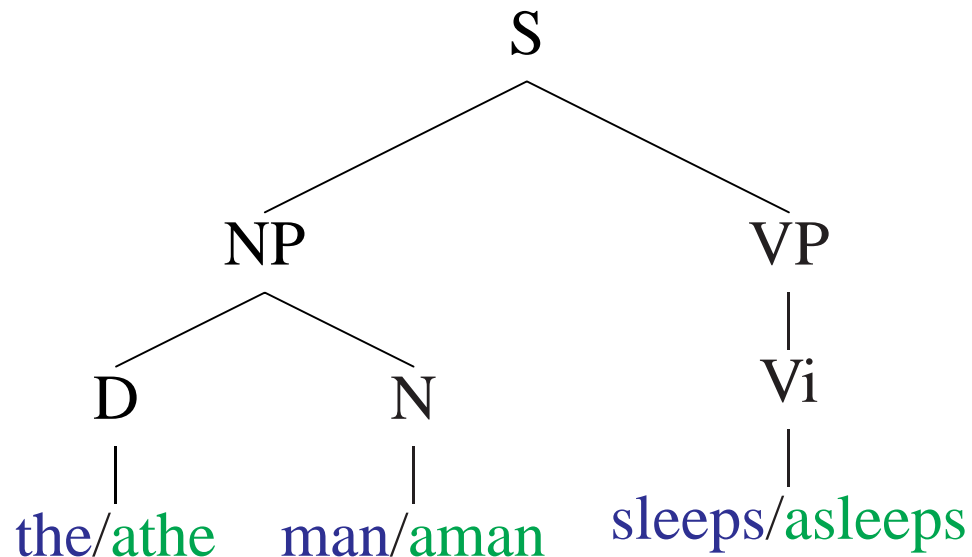
- Probability of a tree with rules $\alpha_i \rightarrow \beta_i$ is $\prod_i P(\alpha_i \rightarrow \beta_i | \alpha_i)$

Transduction PCFGs

- First change to the rules: **lexical** rules generate a pair of words

Vi	⇒	sleeps/ asleeps	1.0
Vt	⇒	saw/ asaw	1.0
NN	⇒	man/ aman	0.7
NN	⇒	woman/ awoman	0.2
NN	⇒	telescope/ atelescope	0.1
DT	⇒	the/ athe	1.0
IN	⇒	with/ awith	0.5
IN	⇒	in/ ain	0.5

Transduction PCFGs



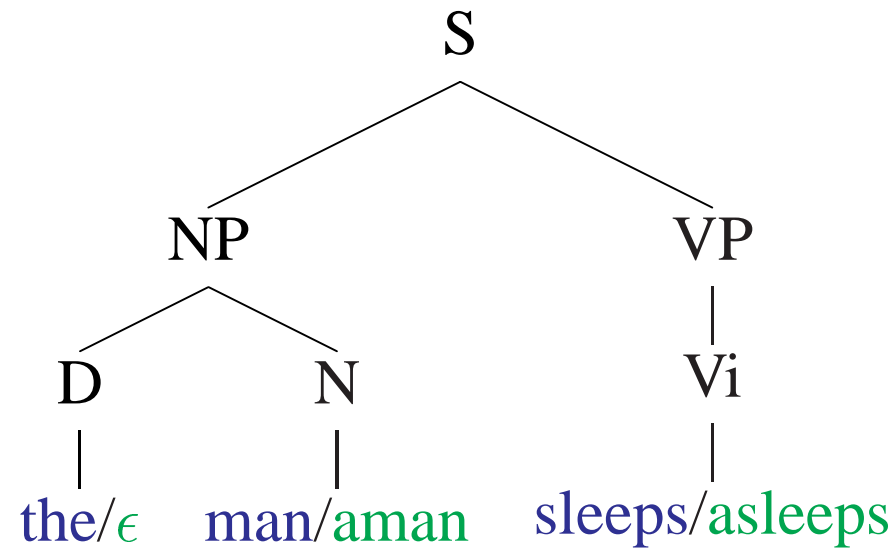
- The modified PCFG gives a distribution over (f, e, T) triples, where e is an English string, f is a French string, and T is a tree

Transduction PCFGs

- Another change: allow empty string ϵ to be generated in either language, e.g.,

DT	\Rightarrow	the/ ϵ	1.0
IN	\Rightarrow	ϵ /awith	0.5

Transduction PCFGs



- Allows strings in the two languages to have different lengths

the man sleeps \Rightarrow aman asleeps

Transduction PCFGs

- Final change: currently formalism does not allow different word orders in the two languages
- Modify the method to allow two types of rules, for example

$$S \Rightarrow [NP \quad VP] \quad 0.7$$

$$S \Rightarrow \langle NP \quad VP \rangle \quad 0.3$$

- Define:
 - E_X is the English string under non-terminal X
e.g., E_{NP} is the English string under the NP
 - F_X is the French string under non-terminal X
- Then for $S \Rightarrow [NP \ VP]$ we define

$$E_S = E_{NP}.E_{VP}$$

$$F_S = F_{NP}.F_{VP}$$

where $.$ is concatenation operation

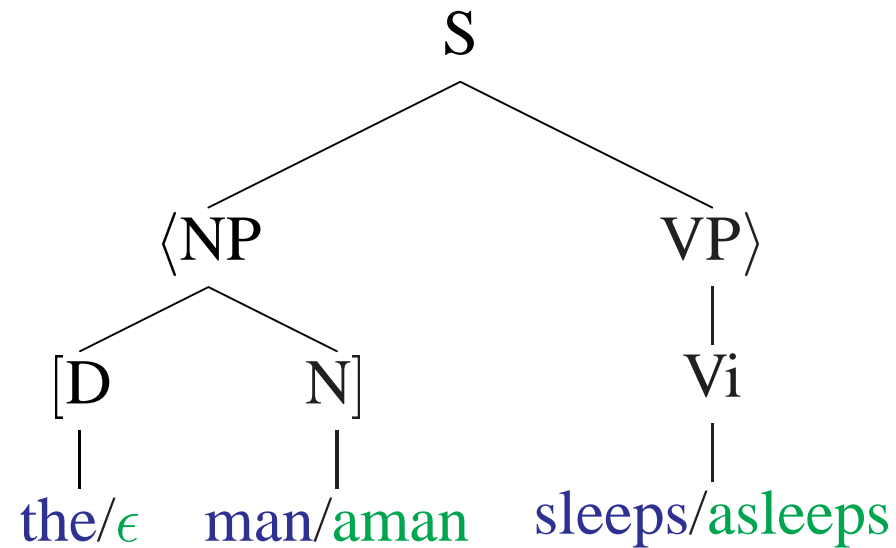
- For $S \Rightarrow \langle NP \ VP \rangle$ we define

$$E_S = E_{NP}.E_{VP}$$

$$F_S = F_{VP}.F_{NP}$$

In the second case, the string order in French is reversed

Transduction PCFGs



- This tree represents the correspondance

the man sleeps \Rightarrow asleeps aman

A Transduction PCFG

S	\Rightarrow	[NP VP]	0.7
S	\Rightarrow	\langle NP VP \rangle	0.3
VP	\Rightarrow	Vi	0.4
VP	\Rightarrow	[Vt NP]	0.01
VP	\Rightarrow	\langle Vt NP \rangle	0.79
VP	\Rightarrow	[VP PP]	0.2
NP	\Rightarrow	[DT NN]	0.55
NP	\Rightarrow	\langle DT NN \rangle	0.15
NP	\Rightarrow	[NP PP]	0.7
PP	\Rightarrow	\langle P NP \rangle	1.0

Vi	⇒	sleeps/ε	0.4
Vi	⇒	sleeps/asleeps	0.6
Vt	⇒	saw/asaw	1.0
NN	⇒	ε/aman	0.7
NN	⇒	woman/awoman	0.2
NN	⇒	telescope/atelescope	0.1
DT	⇒	the/athe	1.0
IN	⇒	with/awith	0.5
IN	⇒	in/ain	0.5

(Wu 1995)

- Dynamic programming algorithms exist for “parsing” a pair of English/French strings (finding most likely tree underlying an English/French pair). Runs in $O(|\mathbf{e}|^3|\mathbf{f}|^3)$ time.
- Training the model: given $(\mathbf{e}_k, \mathbf{f}_k)$ pairs in training data, the model gives

$$P(T, \mathbf{e}_k, \mathbf{f}_k \mid \Theta)$$

where T is a tree, Θ are the parameters. Also gives

$$P(\mathbf{e}_k, \mathbf{f}_k \mid \Theta) = \sum_T P(T, \mathbf{e}_k, \mathbf{f}_k \mid \Theta)$$

Likelihood function is then

$$L(\Theta) = \sum_k \log P(f_k, e_k \mid \Theta) = \sum_k \log \sum_T P(T, f_k, e_k \mid \Theta)$$

Wu gives a dynamic programming implementation for EM

Overview

- Syntax Based Model 1: (Yamada and Knight 2001)
- Syntax Based Model 2: (Wu 1995)
- **A Phrase-Based Model:** (Koehn, Och and Marcu 2003)

Methods that go beyond word-word alignments

A Phrase-Based Model

(Koehn, Och and Marcu 2003)

- Intuition: IBM models have word-word translation
- Intuition: in IBM models each French word is aligned with only one English word
- A new type of model:
align phrases in English with phrases in French

- An example from Koehn and Knight tutorial:

Morgen fliege ich nach Kanada zur Konferenz

Tomorrow I will fly to the conference in Canada

Morgen

fliege

ich

nach Kanada

zur Konferenz

Tomorrow

will fly

I

in Canada

to the conference

Representation as Alignment “Matrix”

	Maria	no	daba	una	bof'	a	la	bruja	verde
Mary	●								
did		●							
not		●							
slap			●	●	●				
the						●	●		
green									●
witch								●	

(Note: “bof’” = “bofetada”)

(Another example from the Koehn and Knight tutorial)

The Issues Involved

- Finding alignment matrices for all English/French pairs in training corpora
- Coming up with a model that incorporates phrases
- Training the model
- Decoding with the model

Finding Alignment Matrices

- Step 1: train model 4 for $P(f | e)$, and come up with most likely alignment for each (e, f) pair
- Step 2: train model 4 for $P(e | f)(!)$ and come up with most likely alignment for each (e, f) pair
- We now have two alignments:
take intersection of the two alignments as a starting point

Intersection of the two alignments:

	Maria	no	daba	una	bof'	a	la	bruja	verde
Mary	●								
did									
not		●							
slap					●				
the							●		
green									●
witch								●	

The intersection of the two alignments was found to be a very reliable starting point

Heuristics for Growing Alignments

- Only explore alignment in **union** of $P(f \mid e)$ and $P(e \mid f)$ alignments
- Add one alignment point at a time
- Only add alignment points which align a word that currently has no alignment
- At first, restrict ourselves to alignment points that are “neighbors” (adjacent or diagonal) of current alignment points
- Later, consider other alignment points

The Issues Involved

- Finding alignment matrices for all English/French pairs in training corpora
- Coming up with a model that incorporates phrases
- Training the model
- Decoding with the model

The Model

- The probability model again models $P(\mathbf{f} \mid \mathbf{e})$
- The steps:
 - Choose a segmentation of \mathbf{e} (all segmentations are equally likely)
 - For each English phrase e , choose a French phrase f with probability

$$\mathbf{T}(f \mid e)$$

for example

$$\mathbf{T}(\text{daba una bofetada} \mid \text{slap})$$

- Choose positions for the French phrases: if start position of the i 'th French phrase is a_i , and end point of $(i - 1)$ 'th French phrase is b_{i-1} , then this has probability

$$\mathbf{R}(a_i - b_{i-1})$$

Training the Model

Simple once we have the alignment matrices!:

- Take maximum-likelihood estimates, e.g.,

$$\mathbf{T}(\text{daba una bofetada} \mid \text{slap}) = \frac{\textit{Count}(\text{daba una bofetada, slap})}{\textit{Count}(\text{slap})}$$

- Take similar estimates for the alignment probabilities

The Issues Involved

- Finding alignment matrices for all English/French pairs in training corpora
- Coming up with a model that incorporates phrases
- Training the model
- Decoding with the model

The Decoding Method

- Goal is to find a high probability English string \mathbf{e} under

$$P(\mathbf{e})P(\mathbf{f}, \mathbf{a} \mid \mathbf{e})$$

where

$$P(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) = \prod_{i=1}^n \mathbf{T}(f_i \mid e_i) \mathbf{R}(a_i - b_{i-1})$$

where f_i and e_i are the n *phrases* in the alignment,
 a_i and b_i are start/end points of the i 'th phrase

The Decoding Method

- A **partial hypothesis** is an English prefix, aligned with some of the French sentence

Maria no daba una bofetada a la bruja verde
| |
Mary did not

- S_m is a **stack** which stores n most likely partial analyses that account for m French words
- At each point, pick a partial hypothesis, and **advance** it by choosing a substring of the French string

Maria no daba una bofetada a la bruja verde

Mary did not



Maria no daba una bofetada a la bruja verde

Mary did not slap

- In this case, we create a new member of the stack S_5