# Cooperative unsupervised training of the part-of-speech taggers in a bidirectional machine translation system[*]

Felipe Sánchez-Martínez, Juan Antonio Pérez-Ortiz, Mikel L. Forcada
Departament de Llenguatges i Sistemes Informàtics
Universitat d'Alacant
E-03071 Alacant, Spain

{fsanchez,japerez,mlf}@dlsi.ua.es

# Contents

- Introduction

- Part-of-speech ambiguities in machine translation

- Part-of-speech tagging with HMM

- Target-language based training of HMM-based taggers

- Cooperative learning of HMM

- Experiments

- Discussion

- Future work

# Introduction

**Part-of-speech (PoS) tagging:** To determine the lexical category or PoS of each word that appears in a text

**Ambiguous word:** Word with more than one possible lexical category (PoS)

|        | **Lemma** | **PoS** |
|--------|-----------|---------|
| *book* | *book*    | noun    |
|        | *book*    | verb    |

Ambiguities are usually solved by looking at the context

# PoS ambiguities in machine translation (I)

**Indirect MT system:** Source language (SL) text is analysed and transformed into an abstract intermediate representation, transformations are applied and, finally, target language (TL) text is generated.

SLAR ↓    TLAR ↓

SL text ⟶ | Analysis | ⟶ | Transformation | ⟶ | Generation | ⟶ TL text

- Analysis module usually includes a PoS tagger

# PoS ambiguities in machine translation (II)

**Mistranslation due to wrong PoS tagging**

- Translation differs from one PoS to another:

| Spanish | PoS | Translation into Catalan |
|---|---|---|
| *para* | preposition | *per a* (for/to) |
| | verb | *para* (stop) |

# PoS ambiguities in machine translation (II)

## Mistranslation due to wrong PoS tagging

- Translation differs from one PoS to another:

| Spanish | PoS | Translation into Catalan |
|---------|-----|--------------------------|
| *para* | preposition | *per a* (for/to) |
| | verb | *para* (stop) |

- Some transformation is applied (or not) for some PoS:

| Spanish | PoS | Translation into Catalan | |
|---------|-----|--------------------------|---|
| *las calles* | *la* (article) | *els carrers* (the streets) | gender |
| | *la* (pronoun) | * *les carrers* (them streets) | ←agreement |

gender
←agreement
rule applied

# PoS tagging with HMM (I)

Use of a hidden Markov model (HMM):

- Adopting a reduced tag set (grouping the finer tags delivered by the morphological analyser)

- Each HMM state corresponds to a different PoS tag

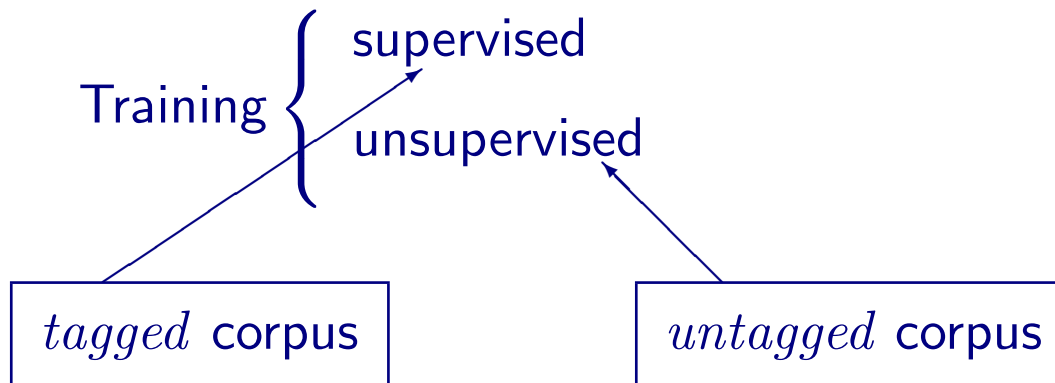- Each input word is replaced by its corresponding ambiguity class

# PoS tagging with HMM (II)

Estimating proper HMM parameters

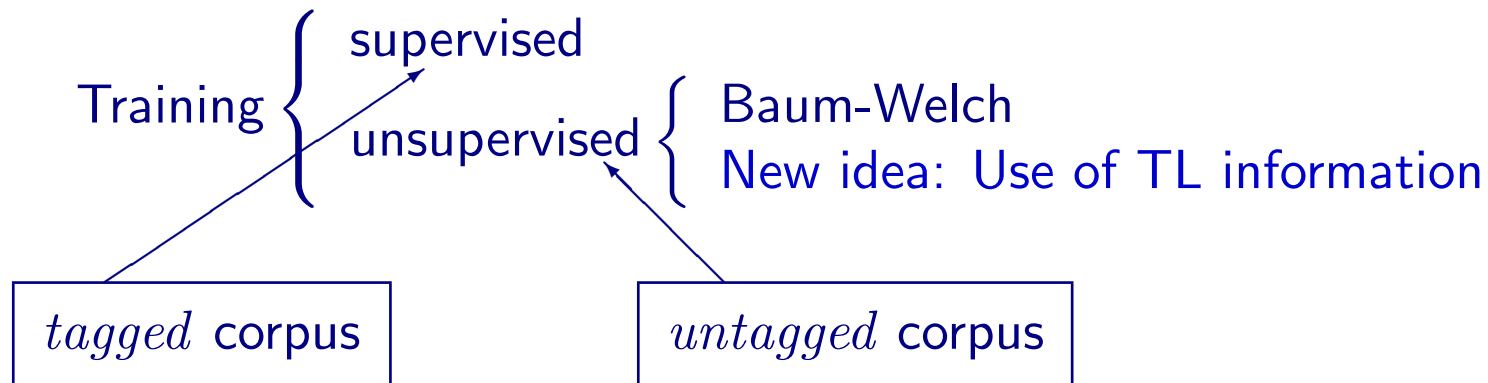$$\text{Training} \begin{cases} \text{supervised} \\ \\ \text{unsupervised} \end{cases}$$

# PoS tagging with HMM (II)

Estimating proper HMM parameters

# PoS tagging with HMM (II)

Estimating proper HMM parameters

Training $\Big\{$ supervised

unsupervised $\Big\{$ Baum-Welch

New idea: Use of TL information

$\boxed{\textit{tagged}\ \text{corpus}}$ $\boxed{\textit{untagged}\ \text{corpus}}$

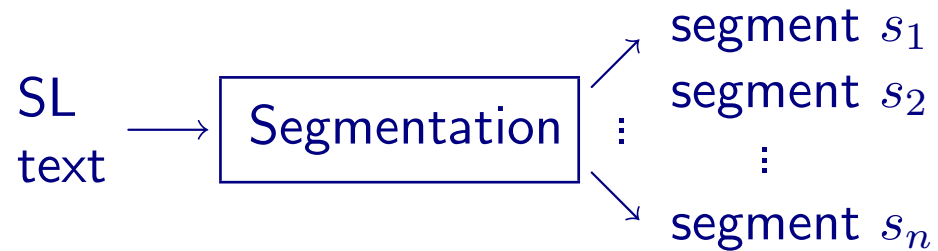# Target-language based training of HMM-based taggers (I)

- Transition probabilities

$$a_{\gamma_i \gamma_j} = \frac{\tilde{n}(\gamma_i \gamma_j)}{\sum_{\gamma_k \in \Gamma} \tilde{n}(\gamma_i \gamma_k)}$$
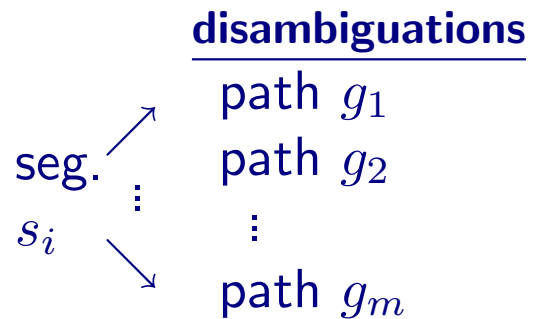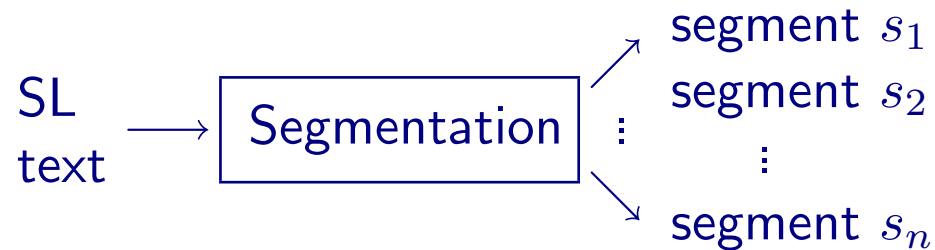
- Emission probabilities

$$b_{\gamma_i \sigma} = \frac{\tilde{n}(\sigma, \gamma_i)}{\sum_{\sigma' : \gamma_i \in \sigma'} \tilde{n}(\sigma', \gamma_i)}$$

# Target-language based training of HMM-based taggers (II)

SL text $\longrightarrow$ [ Segmentation ] $\vdots$   segment $s_1$

segment $s_2$

$\vdots$

segment $s_n$

# Target-language based training of HMM-based taggers (II)

$$\text{SL text} \longrightarrow \boxed{\text{Segmentation}} \;\vdots\; \begin{array}{l} \text{segment } s_1 \\ \text{segment } s_2 \\ \vdots \\ \text{segment } s_n \end{array}$$

$$s_i \text{ seg.} \; \vdots \; \begin{array}{l} \underline{\textbf{disambiguations}} \\ \text{path } g_1 \\ \text{path } g_2 \\ \vdots \\ \text{path } g_m \end{array}$$

# Target-language based training of HMM-based taggers (II)

$$\text{SL text} \longrightarrow \boxed{\text{Segmentation}} \Bigg\langle \begin{array}{l} \text{segment } s_1 \\ \text{segment } s_2 \\ \vdots \\ \text{segment } s_n \end{array}$$

**disambiguations**

**translations**

$$\text{seg. } s_i \Bigg\langle \begin{array}{l} \text{path } g_1 \\ \text{path } g_2 \\ \vdots \\ \text{path } g_m \end{array} \Bigg\rangle \boxed{\text{MT}} \Bigg\langle \begin{array}{l} \tau(g_1, s) \\ \tau(g_2, s) \\ \vdots \\ \tau(g_m, s) \end{array}$$

# Target-language based training of HMM-based taggers (II)

SL
text $\longrightarrow$ | Segmentation | ⋮

segment $s_1$
segment $s_2$
⋮
segment $s_n$

**disambiguations**

path $g_1$
path $g_2$
⋮
path $g_m$

**translations**

$\tau(g_1, s)$
$\tau(g_2, s)$
⋮
$\tau(g_m, s)$

**likelihoods**

$p_{\mathrm{TL}}(\tau(g_1, s))$
$p_{\mathrm{TL}}(\tau(g_2, s))$
⋮
$p_{\mathrm{TL}}(\tau(g_m, s))$

seg.
$s_i$

| MT |

| TL model |

# Target-language based training of HMM-based taggers (II)

$$\text{SL text} \longrightarrow \boxed{\text{Segmentation}} \nearrow \text{segment } s_1$$

segment $s_1$
segment $s_2$
$\vdots$
segment $s_n$

**disambiguations**

path $g_1$
path $g_2$
$\vdots$
path $g_m$

seg.
$s_i$

$\boxed{\text{MT}}$

**translations**

$\tau(g_1, s)$
$\tau(g_2, s)$
$\vdots$
$\tau(g_m, s)$

$\boxed{\begin{array}{c}\text{TL}\\\text{model}\end{array}}$

**likelihoods**

$p_{\mathsf{TL}}(\tau(g_1, s)) \dashrightarrow p(g_1|s)$
$p_{\mathsf{TL}}(\tau(g_2, s)) \dashrightarrow p(g_2|s)$
$\vdots \qquad\qquad \vdots$
$p_{\mathsf{TL}}(\tau(g_m, s)) \dashrightarrow p(g_m|s)$

# Target-language based training of HMM-based taggers (III)

| $s \equiv$ | $y$ | $la$ | $para$ | $si$ | |
|---|---|---|---|---|---|
| | $\{$ CNJ $\}$ | $\left\{ \begin{array}{c} \text{ART} \\ \text{PRN} \end{array} \right\}$ | $\left\{ \begin{array}{c} \text{VB} \\ \text{PR} \end{array} \right\}$ | $\{$ CNJ $\}$ | |
| | | | | | $p(g_i\|s)$ |
| $g_1 \equiv$ | CNJ | ART | PR | CNJ | |
| $\tau(g_1, s) \equiv$ | $i$ (and) | $la$ (the) | $per\ a$ (for/to) | $si$ (if) | 0.0001 |
| $g_2 \equiv$ | CNJ | ART | VB | CNJ | |
| $\tau(g_2, s) \equiv$ | $i$ (and) | $la$ (the) | $para$ (stop) | $si$ (if) | 0.4999 |
| $g_3 \equiv$ | CNJ | PRN | PR | CNJ | |
| $\tau(g_3, s) \equiv$ | $i$ (and) | $la$ (it/her) | $per\ a$ (for/to) | $si$ (if) | 0.0001 |
| $g_4 \equiv$ | CNJ | PRN | VB | CNJ | |
| $\tau(g_4, s) \equiv$ | $i$ (and) | $la$ (it/her) | $para$ (stop) | $si$ (if) | 0.4999 |

# Target-language based training of HMM-based taggers (III)

| $s \equiv$ | $y$ | $la$ | $para$ | $si$ | |
|---|---|---|---|---|---|
| | $\{ \text{ CNJ } \}$ | $\left\{ \begin{array}{c} \text{ART} \\ \text{PRN} \end{array} \right\}$ | $\left\{ \begin{array}{c} \text{VB} \\ \text{PR} \end{array} \right\}$ | $\{ \text{ CNJ } \}$ | |
| | | | | | $p(g_i|s)$ |
| $g_1 \equiv$ | CNJ | ART | PR | CNJ | |
| $\tau(g_1, s) \equiv$ | $i$ (and) | $la$ (the) | $per\ a$ (for/to) | $si$ (if) | 0.0001 |
| $g_2 \equiv$ | CNJ | ART | VB | CNJ | |
| $\tau(g_2, s) \equiv$ | $i$ (and) | $la$ (the) | $para$ (stop) | $si$ (if) | 0.4999 |
| $g_3 \equiv$ | CNJ | PRN | PR | CNJ | |
| $\tau(g_3, s) \equiv$ | $i$ (and) | $la$ (it/her) | $per\ a$ (for/to) | $si$ (if) | 0.0001 |
| $g_4 \equiv$ | CNJ | PRN | VB | CNJ | |
| $\tau(g_4, s) \equiv$ | $i$ (and) | $la$ (it/her) | $para$ (stop) | $si$ (if) | 0.4999 |

**Free ride:** word translated the same way independently of the tag selected

# Target-language based training of HMM-based taggers (IV)

$$p(g_i|s) \propto p(g_i|\tau(g_i, s))\, p_{\mathsf{TL}}(\tau(g_i, s))$$

- $p(g_i|s)$: Probability of path $g_i$ to be the correct disambiguation of segment $s$

- $p_{\mathsf{TL}}(\tau(g_i, s))$: Likelihood of the translation into TL of segment $s$ according to the disambiguation given by path $g_i$

    - Language model based on trigrams of words
    - Hidden Markov model
    - ...

- $p(g_i|\tau(g_i, s))$: Contribution of the disambiguation path $g_i$ to the translation given by $\tau(g_i, s)$

# Cooperative learning of HMM (I)

• Use of the prevoius idea ...

# Cooperative learning of HMM (I)

- Use of the prevoius idea ...

- Bidirectional MT system translating between languages $A$ and $B$

# Cooperative learning of HMM (I)

- Use of the prevoius idea ...

- Bidirectional MT system translating between languages $A$ and $B$

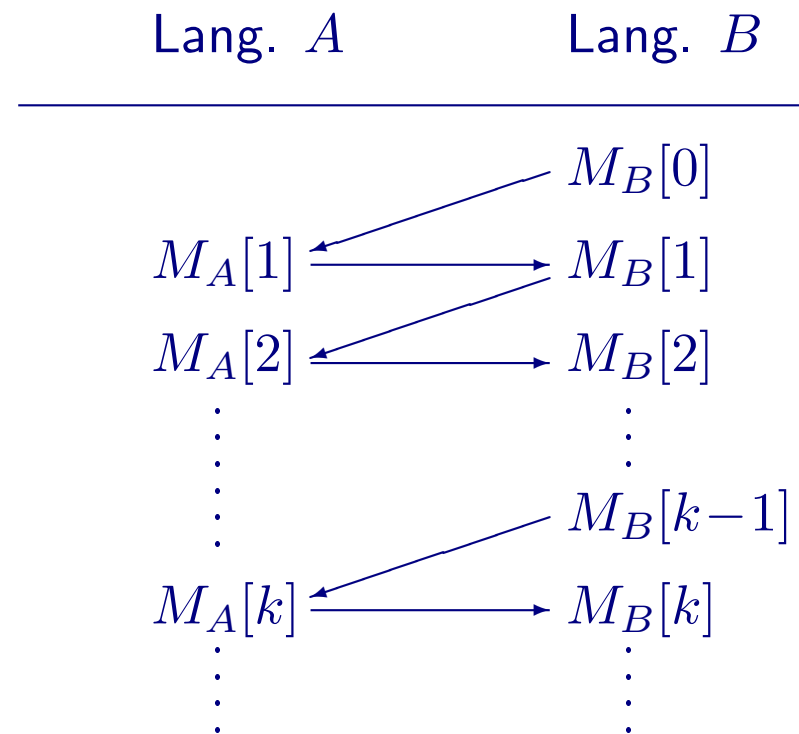- Morphological generation is not done when performing translations

# Cooperative learning of HMM (I)

- Use of the prevoius idea ...

- Bidirectional MT system translating between languages $A$ and $B$

- Morphological generation is not done when performing translations

- Before morphological generation we have a sequence of lexical categories (tags) in the TL

# Cooperative learning of HMM (I)

- Use of the prevoius idea ...

- Bidirectional MT system translating between languages $A$ and $B$

- Morphological generation is not done when performing translations

- Before morphological generation we have a sequence of lexical categories (tags) in the TL

- Use of such a TL model based on tags: HMM as a TL model

# Cooperative learning of HMM (II)

$$\begin{array}{cc}
\text{Lang. } A & \text{Lang. } B \\
\hline
& M_B[0] \\
M_A[1] & M_B[1] \\
M_A[2] & M_B[2] \\
\vdots & \vdots \\
& M_B[k{-}1] \\
M_A[k] & M_B[k] \\
\vdots & \vdots
\end{array}$$

# Experiments

- We used the Spanish↔Catalan MT system interNOSTRUM
  (`www.internostrum.com`)

  Language $A$: Catalan
  Language $B$: Spanish

# Experiments

- We used the Spanish↔Catalan MT system interNOSTRUM (`www.internostrum.com`)

  Language $A$: Catalan
  Language $B$: Spanish

- Use of various corpus sizes and three different corpora for each size

# Experiments

- We used the Spanish↔Catalan MT system interNOSTRUM (`www.internostrum.com`)

  Language $A$: Catalan
  Language $B$: Spanish

- Use of various corpus sizes and three different corpora for each size

- Evaluation with an independent corpus for each language:

  – PoS error rate with hand-tagged corpus
  – Translation error rate with human-corrected translations

# Results

Proof of two different initial models $M_B[0]$:

- "Good" HMM: Trained from $1\,000\,000$-word untagged SL corpus with the Baum-Welch algorithm (PoS error rate: $34.2\%$)

- "Bad" HMM: Equiprobable transition and emission probabilities (PoS error rate: $76.5\%$)

# Results

Proof of two different initial models $M_B[0]$:

- "Good" HMM: Trained from $1\,000\,000$-word untagged SL corpus with the Baum-Welch algorithm (PoS error rate: $34.2\%$)

- "Bad" HMM: Equiprobable transition and emission probabilities (PoS error rate: $76.5\%$)

|  | Avg. PoS error | | Avg. translation error | | Avg. It. |
|---|---|---|---|---|---|
|  | Spanish | Catalan | Spanish | Catalan |  |
| "good start" → | 24.9% | 27.5% | 6.2% | 6.7% | 2 |
| "bad start" → | 25.9% | 26.4% | 6.1% | 6.8% | 5 |
| Baum-Welch → | 31.7% | 37.8% | 8.4% | 13.6% | 14 |
| supervised → | 10.4% | 16.5% | 2.6% | 3.0% |  |

# Discussion

- PoS error and translation error rates lie between those produced by supervised and unsupervised methods

# Discussion

- PoS error and translation error rates lie between those produced by supervised and unsupervised methods

- There is no need for good initial information to achieve good results

# Discussion

- PoS error and translation error rates lie between those produced by supervised and unsupervised methods

- There is no need for good initial information to achieve good results

- The method described needs a relatively small amount of words compare with common corpus sizes used with the Baum-Welch algorithm

# Discussion

- PoS error and translation error rates lie between those produced by supervised and unsupervised methods

- There is no need for good initial information to achieve good results

- The method described needs a relatively small amount of words compare with common corpus sizes used with the Baum-Welch algorithm

- The training method produces PoS taggers tuned not only with SL texts, but also with TL texts and the underlying MT system

# Future work

- Research on better estimates for $p(g_i|\tau(g_i, s))$

  - Estimate the HMM parameters iteratively
    Use the parameters of the previous iteration to estimate $p(g_i|\tau(g_i, s))$

# Future work

- Research on better estimates for $p(g_i|\tau(g_i, s))$

  - Estimate the HMM parameters iteratively
    Use the parameters of the previous iteration to estimate $p(g_i|\tau(g_i, s))$

- Time complexity reduction

  - Use of a $k$-best Viterbi algorithm with the current parameters to calculate approximate likelihood and translate only the $k$ most promising paths

# Future work

- Research on better estimates for $p(g_i|\tau(g_i, s))$

  - Estimate the HMM parameters iteratively
    Use the parameters of the previous iteration to estimate $p(g_i|\tau(g_i, s))$

- Time complexity reduction

  - Use of a $k$-best Viterbi algorithm with the current parameters to calculate approximate likelihood and translate only the $k$ most promising paths

- Better formalization

  - Different disambiguation paths from different segments can produce the same translation

# Thank you very much for your attention !!