Machine Learning applied to Natural Language Processing

Lluís Màrquez

ALT Research Group Qatar Computing Research Institute

Advanced Methods for Corpus Analysis LCT – European Masters Program in Language and Communication Technologies Donostia, June 22-24, 2015

Talk Overview



2 Semantic Role Labeling — A Running Example

3 Conclusions



Natural Language Processing Applications

• Typical NLP applications:

- Machine Translation
- (CL) Information Retrieval and document managment
- Information Extraction
- Modern Question Answering (e.g., Watson)
- Machine Reading
- Document Summarization (multidocument, multilingual)
- Dialog Systems
- Different levels of linguistic knowledge and comprehension are required
- They need to resolve a number of basic subproblems



Natural Language Processing Applications

• Typical NLP applications:

- Machine Translation
- (CL) Information Retrieval and document managment
- Information Extraction
- Modern Question Answering (e.g., Watson)
- Machine Reading
- Document Summarization (multidocument, multilingual)
- Dialog Systems
- Different levels of linguistic knowledge and comprehension are required
- They need to resolve a number of basic subproblems



Simple Idea:

- Mapping from an input to an output structure
 - The input structure is typically a sequence of words, which might be enriched with some linguistic information
 - Output structures are sequences, trees, graphs, etc.



Part-of-Speech Tagging

The San Francisco Examiner issued a special edition around noon yesterday that was filled entirely with earthquake new and information.



Part-of-Speech Tagging

The_DT San_NNP Francisco_NNP Examiner_NNP issued_VBD a_DT special_JJ edition_NN around_IN noon_NN yesterday_NN that_WDT was_VBD filled_VBN entirely_RB with_IN earthquake_NN news_NN and_CC information_NN ._.

POS tagging is a pure sequential labeling problem

(sequential learning paradigm)

But... are really words ambiguous with respect to POS?



Part-of-Speech Tagging

The_DT San_NNP Francisco_NNP Examiner_NNP issued_VBD a_DT special_JJ edition_NN around_IN noon_NN yesterday_NN that_WDT was_VBD filled_VBN entirely_RB with_IN earthquake_NN news_NN and_CC information_NN ._.

POS tagging is a pure sequential labeling problem

(sequential learning paradigm)

But... are really words ambiguous with respect to POS?



Part-of-Speech Tagging

The_DT San_NNP Francisco_NNP Examiner_NNP issued_VBD a_DT special_JJ edition_NN around_IN noon_NN yesterday_NN that_WDT was_VBD filled_VBN entirely_RB with_IN earthquake_NN news_NN and_CC information_NN ._.

But... are really words ambiguous with respect to POS?

YES! Let's take a look at a free on-line demo: **FreeLing** http://nlp.lsi.upc.edu/freeling/demo/demo.php



Syntactic Analysis (Constituency parsing)

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



Syntactic Analysis (Constituency parsing)

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



Dependency Parsing





Shallow Parsing (Chunking)

He reckons the current account deficit will narrow to only 1.8 billion in September.



Shallow Parsing (Chunking)

 $\label{eq:product} \begin{array}{l} \left[\text{NP He} \right] \left[\text{VP reckons} \right] \left[\text{NP the current account deficit} \right] \left[\text{VP will narrow} \right] \left[\text{PP o} \right] \left[\text{NP only 1.8 billion} \right] \left[\text{PP in} \right] \left[\text{NP September} \right] . \end{array}$

Chunking is a sequential phrase recognition task

It can be seen as a sequential labeling problem (B-I-O encoding)

He_B-NP reckons_B-VP the_B-NP current_I-NP account_I-NP deficit_I-NP will_B-VP narrow_I-VP to_B-PP only_B-NP 1.8_I-NP billion_I-NP in_B-PP September_B-NP ._O

this is simple and usually effective



Shallow Parsing (Chunking)

 $\label{eq:product} \begin{array}{l} \left[\text{NP He} \right] \left[\text{VP reckons} \right] \left[\text{NP the current account deficit} \right] \left[\text{VP will narrow} \right] \left[\text{PP o} \right] \left[\text{NP only 1.8 billion} \right] \left[\text{PP in} \right] \left[\text{NP September} \right] . \end{array}$

Chunking is a sequential phrase recognition task

It can be seen as a sequential labeling problem (B-I-O encoding)

He_B-NP reckons_B-VP the_B-NP current_I-NP account_I-NP deficit_I-NP will_B-VP narrow_I-VP to_B-PP only_B-NP 1.8_I-NP billion_I-NP in_B-PP September_B-NP ._O

this is simple and usually effective



Natural Language Processing Problems₍₄₎

Clause splitting (partial parsing)

The deregulation of railroads and trucking companies that began in 1980 enabled shippers to bargain for transportation.



Natural Language Processing Problems₍₄₎

Clause splitting (partial parsing)

(s The deregulation of railroads and trucking companies (SBAR that (s began in 1980)) enabled (s shippers to bargain for transportation) .)



Natural Language Processing Problems₍₄₎

Clause splitting (partial parsing)

```
(S The deregulation of railroads and trucking companies
        (SBAR that
            (S began in 1980 ))
            enabled
            (S shippers to bargain for transportation)
. )
```

Clauses may embed: they form a hierarchy

Clause splitting is a hierarchical prhase recognition problem

Not a good idea to treat it as a sequential problem...



Semantic Role Labeling (shallow semantic parsing)

He wouldn't accept anything of value from those he was writing about.



Semantic Role Labeling (shallow semantic parsing)

 $\begin{bmatrix} A_0 & \text{He} \end{bmatrix} \begin{bmatrix} AM-MOD & \text{would} \end{bmatrix} \begin{bmatrix} AM-NEG & n't \end{bmatrix} \begin{bmatrix} V & \text{accept} \end{bmatrix} \begin{bmatrix} A_1 & \text{anything of value} \end{bmatrix} \text{ from } \begin{bmatrix} A_2 & \text{those he was writing about} \end{bmatrix} .$

Roles for the predicate accept (PropBank Frames scheme):

V: verb; A₀: acceptor; A₁: thing accepted; A₂: accepted-from; A₃: attribute; AM-MOD: modal; AM-NEG: negation;



Motivation

Natural Language Processing Problems(5)





Named Entity Extraction ("semantic chunking")

Wolff, currently a journalist in Argentina, played with Del Bosque in the final years of the seventies in Real Madrid.



Named Entity Extraction ("semantic chunking")

 $[\ensuremath{\text{PER}}$ Wolff] , currently a journalist in [LOC Argentina] , played with [PER Del Bosque] in the final years of the seventies in [ORG Real Madrid] .

- Named Entities may be embedded
- NE tracing: variants and co-reference resolution
- Relations between entities: event extraction



Named Entities, relations, events, etc. (example from the ACE corpus)

LOS ANGELES, April 18 (AFP)

Best-selling novelis and "Jurassic Park" creator Michael Crichton has agreed to pay hisfourth wife 31 million dollars as part of their divorce settlement, court documents showed Friday. Crichton 60 is one of the world's wealthiest authors, and has had 12 of his novels made into major Hollywood movies. The writer will retain the rights to his books and films, although he has agreed to split a raft of other possessions with Anne Marie his wife of 13 years according to documents filed in Los Angeles Superior Court.





Discourse Parsing





Recall the take away message:

- Mapping from an input to an output structure
 - The input structure is typically a sequence of words, which might be enriched with some linguistic information
 - Output structures are sequences, trees, graphs, etc.
- Machine Learning and Search (inference) are in between



NLP Meets Machine Learning

- 1980's resurgence of the empirical paradigm for NLP
- 1990's massive application of Machine Learning techniques
- Important factor (among others):

• Ambiguity resolution can be directly casted as classification

- NLP community learnt very well how to model and learn local decisions
- Note 1: There is a big gap between classification and structure learning. Pure classification tasks don't really exist!
- Note 2: Search is strongly related to the generation of the output structure (*decoding*, *inference*, etc.)



Why applying Machine Learning?

- Low cost development of linguistic processors
- Language (quasi)independence: reusability
- Ability of acquiring/discovering knowledge from very large datasets
- Assist manual development of linguistic resources



On-line Demos in the Web

- FreeLing. Universitat Politècnica de Catalunya. Basic syntactic processing. Catalan, Spanish, English and others. http://nlp.lsi.upc.edu/freeling/demo/demo.php
- **CCG tools**. University of Illinois at Urbana-Champaign. Multiple processors and applications. English. http://cogcomp.cs.illinois.edu/page/demos/



Talk Overview

1 Motivation

Semantic Role Labeling — A Running Example

- The Statistical Approach to SRL
- Examples of SRL Systems
- Feature Engineering
- Semantic Features for SRL
- An Arc-factored Model for Joint Syntactic-SRL Parsing
- Not Addressed in this Course

3 Conclusions



Semantic Role Labeling

SRL ^{def} = identify the *arguments* of a given proposition and assign them semantic labels describing the *roles* they play in the predicate (i.e., recognize predicate argument structures)



IE point of view

SRL ^{def} = detecting basic event structures such as *who* did *what* to *whom, when* and *where*

[The luxury auto maker]_{AGENT} [last year]_{TEMP} sold_P [1,214 cars]_{OBJECT} [in the U.S.]_{LOCATIVE}



IE point of view

SRL $\stackrel{def}{=}$ detecting basic event structures such as *who* did *what* to *whom, when* and *where*

[The luxury auto maker]_{AGENT} [last year]_{TEMP} sold_P [1,214 cars]_{OBJECT} [in the U.S.]_{LOCATIVE}





Example from (Yih & Toutanova, 2006)



Syntactic variations



- Scott was hit by Kristina yesterday with a baseball
- Yesterday, Scott was hit with a baseball by Kristina
- Yesterday Scott was hit by Kristina with a baseball
- Kristina hit Scott with a baseball yesterday
- \Rightarrow All of them share the same semantic representation:

hit(Kristina,Scott,yesterday,with a baseball)

Example from (Yih & Toutanova, 2006)



Structural view

Mapping from input to output structures:

- Input is text (enriched with morpho-syntactic information)
- Output is a sequence of labeled arguments
- Sequential segmenting/labeling problem

" Mr. Smith sent the report to me this morning . '

[Mr. Smith]_{AGENT} sent [the report]_{OBJ} [to me]_{RECIP} [this morning]_{TMP}.

 $Mr_{B-AGENT}$ Smith₁ sent the_{B-OBJ} report₁ to_{B-RECIP} me₁ this_{B-TMP} morning₁ ._O
Structural view

Mapping from input to output structures:

- Input is *text* (enriched with morpho-syntactic information)
- Output is a sequence of labeled arguments

• Sequential segmenting/labeling problem

" Mr. Smith sent the report to me this morning . '

 $[Mr. Smith]_{AGENT} \text{ sent [the report]}_{OBJ} [to me]_{RECIP} [this morning]_{TMP} .$

 $Mr_{B-AGENT}$ Smith₁ sent the_{B-OBJ} report₁ to_{B-RECIP} me₁ this_{B-TMP} morning₁ ._O

Structural view

Mapping from input to output structures:

- Input is *text* (enriched with morpho-syntactic information)
- Output is a sequence of labeled arguments
- Sequential segmenting/labeling problem

" Mr. Smith sent the report to me this morning . '

[Mr. Smith]_{AGENT} sent [the report]_{OBJ} [to me]_{RECIP} [this morning]_{TMP}.

 $Mr_{B-AGENT}$ Smith₁ sent the_{B-OBJ} report₁ to_{B-RECIP} me₁ this_{B-TMP} morning₁ .0



Output is a hierarchy of labeled arguments





Output is a hierarchy of labeled arguments



Linguistic nature of the problem

• Argument identification is strongly related to syntax



• Role labeling is a semantic task

e.g., selectional preferences could play an important role)



Linguistic nature of the problem

• Argument identification is strongly related to syntax



Role labeling is a semantic task

(e.g., selectional preferences could play an important role)



SRL Systems Available

- ASSERT (Automatic Statistical SEmantic Role Tagger) http://cemantix.org/assert.html
- UIUC system demo

http://l2r.cs.uiuc.edu/~cogcomp/srl-demo.php

- SwiRL: state-of-the-art system from CoNLL-2005 http://www.surdeanu.name/mihai
- Shalmaneser: FrameNet-based system from SALSA project http://www.coli.uni-saarland.de/projects/salsa/shal/
- SEMAFOR: Probabilistic Frame(Net)-Semantic Parser http://www.ark.cs.cmu.edu/SEMAFOR/
- Brutus: A CCG-based Semantic Role Labeler http://www.ling.ohio-state.edu/~boxwell/software/brutus.html



- (English) PropBank http://verbs.colorado.edu/~mpalmer/projects/ace.html
- FrameNet

http://framenet.icsi.berkeley.edu

- Korean PropBank http://www.ldc.upenn.edu/
- Chinese PropBank http://verbs.colorado.edu/chinese/cpb/
- AnCora corpus: Spanish and Catalan http://http://clic.ub.edu/ancora/
- Prague Dependency Treebank: Czech http://ufal.mff.cuni.cz/pdt2.0/
- Penn Arabic TreeBank: Arabic http://www.ircs.upenn.edu/arabic/



• (English) PropBank

http://verbs.colorado.edu/~mpalmer/projects/ace.html

FrameNet

http://framenet.icsi.berkeley.edu

- Korean PropBank http://www.ldc.upenn.edu/
- Chinese PropBank http://verbs.colorado.edu/chinese/cpb/
- AnCora corpus: Spanish and Catalan http://http://clic.ub.edu/ancora/
- Prague Dependency Treebank: Czech http://ufal.mff.cuni.cz/pdt2.0/
- Penn Arabic TreeBank: Arabic http://www.ircs.upenn.edu/arabic/



PropBank

- Syntax-based approach: explaining the varied expression of verb arguments within syntactic positions
- Annotation of all verbal predicates in WSJ (Penn Treebank)
- http://verbs.colorado.edu/~mpalmer/projects/ace.html
- Add a semantic layer to the Syntactic Trees



PropBank

- Syntax-based approach: explaining the varied expression of verb arguments within syntactic positions
- Annotation of all verbal predicates in WSJ (Penn Treebank)
- http://verbs.colorado.edu/~mpalmer/projects/ace.html
- Add a semantic layer to the Syntactic Trees



PropBank

- Theory neutral numbered core roles (Arg0, Arg1, etc.)
 - ⇒ Interpretation of roles: verb-specific framesets
 - ⇒ Arg0 and Arg1 usually correspond to prototypical Agent and Patient/Theme roles. Other arguments do not consistently generalize across verbs
 - ⇒ Different senses have different framesets
 - ⇒ Syntactic alternations that preserve meaning are kept toghether in a single frameset
- Closed set of 13 general labels for Adjuncts (e.g., Temporal, Manner, Location, etc.)



PropBank

- Theory neutral numbered core roles (Arg0, Arg1, etc.)
 - \Rightarrow Interpretation of roles: verb-specific framesets
 - ⇒ Arg0 and Arg1 usually correspond to prototypical Agent and Patient/Theme roles. Other arguments do not consistently generalize across verbs
 - ⇒ Different senses have different framesets
 - ⇒ Syntactic alternations that preserve meaning are kept toghether in a single frameset
- Closed set of 13 general labels for Adjuncts (e.g., Temporal, Manner, Location, etc.)



PropBank

- Theory neutral numbered core roles (Arg0, Arg1, etc.)
 - \Rightarrow Interpretation of roles: verb-specific framesets
 - ⇒ Arg0 and Arg1 usually correspond to prototypical Agent and Patient/Theme roles. Other arguments do not consistently generalize across verbs
 - ⇒ Different senses have different framesets
 - ⇒ Syntactic alternations that preserve meaning are kept toghether in a single frameset
- Closed set of 13 general labels for Adjuncts (e.g., Temporal, Manner, Location, etc.)



PropBank: Frame files

(Palmer et al., <u>2005)</u>

- sell.01: commerce: seller Arg0="seller" (agent); Arg1="thing sold" (theme); Arg2="buyer" (recipient); Arg3="price paid"; Arg4="benefactive" [Al Brownstein]_{Arg0} sold [it]_{Arg1} [for \$60 a bottle]_{Arg3}
- sell.02: give up Arg0="entity selling out" [John]_{Arg0} sold out
- sell.03: sell until none is/are left Arg0="seller"; Arg1="thing sold"; ...

[The new Harry Potter]_{Arg1} sold out [within 20 minutes]_{ArgM-TMP}

Jatar Computing Research Institute Media of Oper Newtonia Milander of Oper

PropBank: Frame files

 sell.01: commerce: seller Arg0= "seller" (agent); Arg1= "thing sold" (theme); Arg2= "buyer" (recipient); Arg3= "price paid"; Arg4= "benefactive" [Al Brownstein]_{Arg0} sold [it]_{Arg1} [for \$60 a bottle]_{Arg3}

 sell.02: give up Arg0="entity selling out"

[John]_{Arg0} sold out

 sell.03: sell until none is/are left Arg0= "seller"; Arg1= "thing sold"; ...

[The new Harry Potter]_{Arg1} sold out [within 20 minutes]_{ArgM-TMP}

Applications

Examples of applications of SRL (I)

- Information Extraction (Surdeanu et al., 2003)
- Question & Answering (Narayanan and Harabagiu, 2004; Frank et al., 2007; Shen and Lapata, 2007)
- Automatic Summarization (Melli et al., 2005)
- Coreference Resolution (Ponzetto and Strube, 2006)
- Text Categorization (Person et al., 2010)
- Opinion Expression Detection (Johansson and Moschitti, 2010)



Applications

Examples of applications of SRL (II)

- Machine Translation Evaluation (Giménez and Màrquez, 2007)
- Machine Translation (Boas, 2002; Wu and Fung, 2009a;2009b)
- Textual Entailment

(Tatu & Moldovan, 2005; Burchardt et al., 2007)

- Modeling Early Language Acquisition (Connor et al., 2008;2009)
- Pictorial Communication Systems (Goldberg, et al., 2008)



Empirical Evaluation of SRL Systems

Evaluation Exercises

- Up to 10 evaluation exercises in the last 7 years
 - $\Rightarrow CoNLL-2004/2005 shared tasks$ (Carreras & Màrquez, 2004; 2005)
 - \Rightarrow Senseval–3 (Litkowski, 2004)
 - ⇒ SemEval-2007 (Pradhan et al., 2007; Màrquez et al., 2007) (Baker et al., 2007; Litkowski & Hargraves, 2007)
 - \Rightarrow CoNLL-2008 shared task (Surdeanu et al., 2008)
 - \Rightarrow CoNLL-2009 shared task (Hajič et al., 2009)
 - \Rightarrow SemEval-2010 (Ruppenhofer et al., 2010)



Talk Overview

Motivation

Semantic Role Labeling — A Running Example

- The Statistical Approach to SRL
- Examples of SRL Systems
- Feature Engineering
- Semantic Features for SRL
- An Arc-factored Model for Joint Syntactic-SRL Parsing
- Not Addressed in this Course

3 Conclusions



Step 1: Select argument candidates

- Given a sentence and a designated predicate
- Parse the sentence
- Identify candidates in tree constituents (filtering/pruning)
 - ⇒ Simple heuristic rules can be used, which maintain a high recall (Xue & Palmer, 2004)

Key point: 95% of semantic arguments coincide with unique syntactic constituents in the gold parse tree (PropBank)
 ⇒ Matching is still ~90% when using automatic parsers



Step 1: Select argument candidates

- Given a sentence and a designated predicate
- Parse the sentence
- Identify candidates in tree constituents (filtering/pruning)
 - ⇒ Simple heuristic rules can be used, which maintain a high recall (Xue & Palmer, 2004)
- Key point: 95% of semantic arguments coincide with unique syntactic constituents in the gold parse tree (PropBank)
 - $\Rightarrow\,$ Matching is still ${\sim}90\%$ when using automatic parsers



Step 2: Local scoring of candidates

- Apply classifiers to assign confidence scores to argument candidates (all labels + 'non-argument')
- Candidates are treated independently of each other
- Identification and Classification may be performed separately
 Computational reasons but also modularity in feature engineering
- Many ML paradigms have been used: not big differences
- Features are more important



Step 2: Local scoring of candidates

- Apply classifiers to assign confidence scores to argument candidates (all labels + 'non-argument')
- Candidates are treated independently of each other
- Identification and Classification may be performed separately
 - ⇒ Computational reasons but also modularity in feature engineering
- Many ML paradigms have been used: not big differences
- Features are more important



Step 2: Local scoring of candidates

- Apply classifiers to assign confidence scores to argument candidates (all labels + 'non-argument')
- Candidates are treated independently of each other
- Identification and Classification may be performed separately
 - ⇒ Computational reasons but also modularity in feature engineering
- Many ML paradigms have been used: not big differences
- Features are more important



SRL Architecture: Steps 1 + 2

Scotty said the same words more loudly










































SRL Architecture: Motivating next step (joint scoring)



•••





SRL Architecture: Motivating next step (joint scoring)



•••





SRL Architecture: Step by Step

Step 3: Joint scoring — Paradigmatic examples

- Combine local predictions through ILP to find the best solution according to structural and linguistic constraints (Koomen et al., 2005; Punyakanok et al., 2008)
 -learning +features +search
- Re-ranking of several candidate solutions (Haghighi et al., 2005; Toutanova et al., 2008)
 +learning +features -search
- Global search integrating joint scoring: Tree CRFs

(Cohn & Blunsom, 2005)

+learning +/-features +/-search

SRL Architecture: Step by Step

Step 3: Joint scoring — Paradigmatic examples

- Combine local predictions through ILP to find the best solution according to structural and linguistic constraints (Koomen et al., 2005; Punyakanok et al., 2008)
 -learning +features +search
- Re-ranking of several candidate solutions

```
(Haghighi et al., 2005; Toutanova et al., 2008)
```

+learning +features -search

• Global search integrating joint scoring: Tree CRFs

```
(Cohn & Blunsom, 2005)
```

+learning +/-features +/-search

Semantic Role Labeling — A Running Example

SRL Architecture: Step by Step

Step 4: Post-processing

 Application of a set of heuristic rules to: Correct frequent errors
 Enforce consistency in the solution



Detour to Machine Learning Concepts

What do we need from ML so far?

- Estimate functions to predict the local scores
 - Supervised machine learning for classification
 - Decision Trees, AdaBoost, MaxEnt, Perceptron, SVMs
- Mechanisms to implement a joint inference process (later...)



Talk Overview

Motivation

Semantic Role Labeling — A Running Example

• The Statistical Approach to SRL

Examples of SRL Systems

- Feature Engineering
- Semantic Features for SRL
- An Arc-factored Model for Joint Syntactic-SRL Parsing
- Not Addressed in this Course

3 Conclusions



Semantic Role Labeling — A Running Example

Examples of SRL systems

- Generalized inference with local classifers and constraints ILP approach (Punyakanok et al., 2008)
- Joint System based on Reranking (Toutanova et al., 2008)
- SRL as sequential labeling (Marquez et al., 2005)



Architecture

- Identify argument candidates
 - Pruning (Xue & Palmer, 2004)
 - Argument identification: binary classification (using SNoW)

Classify argument candidates

Argument Classifier: multi-class classification (SNoW)

- Use the estimated probability distribution given by the argument classifier
- Use structural and linguistic constraints
- Infer the optimal global output



Architecture

- Identify argument candidates
 - Pruning (Xue & Palmer, 2004)
 - Argument identification: binary classification (using SNoW)
- Classify argument candidates
 - Argument Classifier: multi-class classification (SNoW)

- Use the estimated probability distribution given by the argument classifier
- Use structural and linguistic constraints
- Infer the optimal global output



Architecture

- Identify argument candidates
 - Pruning (Xue & Palmer, 2004)
 - Argument identification: binary classification (using SNoW)
- Classify argument candidates
 - Argument Classifier: multi-class classification (SNoW)

- Use the estimated probability distribution given by the argument classifier
- Use structural and linguistic constraints
- Infer the optimal global output



- The output of the argument classifier often violates some constraints, especially when the sentence is long
- Finding the best legitimate output is formalized as an optimization problem and solved via Integer Linear Programming (Roth & Yih, 2004)
- Input formed by:
 - The probability estimation (by the argument classifier)
 - Structural and linguistic constraints
- Allows incorporating expressive constraints (non-sequential) on the variables (the arguments types)



- The output of the argument classifier often violates some constraints, especially when the sentence is long
- Finding the best legitimate output is formalized as an optimization problem and solved via Integer Linear Programming (Roth & Yih, 2004)
- Input formed by:
 - The probability estimation (by the argument classifier)
 - Structural and linguistic constraints
- Allows incorporating expressive constraints (non-sequential) on the variables (the arguments types)



- The output of the argument classifier often violates some constraints, especially when the sentence is long
- Finding the best legitimate output is formalized as an optimization problem and solved via Integer Linear Programming (Roth & Yih, 2004)
- Input formed by:
 - The probability estimation (by the argument classifier)
 - Structural and linguistic constraints
- Allows incorporating expressive constraints (non-sequential) on the variables (the arguments types)



- The output of the argument classifier often violates some constraints, especially when the sentence is long
- Finding the best legitimate output is formalized as an optimization problem and solved via Integer Linear Programming (Roth & Yih, 2004)
- Input formed by:
 - The probability estimation (by the argument classifier)
 - Structural and linguistic constraints
- Allows incorporating expressive constraints (non-sequential) on the variables (the arguments types)



Integer Linear Programming Inference

- For each candidate argument a_i (1 ≤ i ≤ n),
 Set up a Boolean variable: a_{i,t} indicating whether a_i is classified as argument type t
- Goal is to maximize: $\sum_{i} \text{score}(a_i = t) \cdot a_{i,t}$ Subject to the (linear) constraints
- If $score(a_i = t) = P(a_i = t)$, the objective is to find the assignment that maximizes the expected number of arguments that are correct and satisfies the constraints



Constraints: examples

- No duplicate argument classes: $\sum_{i=1}^{n} a_{i,Arg0} \leq 1$
- On discontinuous arguments (C-ARG) $\forall j (1 \leq j \leq n), \sum_{i=1}^{j-1} a_{i,Arg0} \geq a_{j,C-Arg0}$
- On reference arguments (R-ARG) $\forall j (1 \le j \le n), \sum_{i \ne j} a_{i,Arg0} \ge a_{j,R-Arg0}$
- Many other possible constraints:
 - Unique labels
 - No overlapping or embedding
 - Relations between number of arguments; order constraints
 - If verb is of type A, no argument of type B
- ILP inference can be used to combine different SRL systems

Constraints: examples

- No duplicate argument classes: $\sum_{i=1}^{n} a_{i,Arg0} \leq 1$
- On discontinuous arguments (C-ARG) $\forall j (1 \le j \le n), \sum_{i=1}^{j-1} a_{i,Arg0} \ge a_{j,C-Arg0}$
- On reference arguments (R-ARG) $\forall j (1 \le j \le n), \sum_{i \ne j} a_{i,Arg0} \ge a_{j,R-Arg0}$
- Many other possible constraints:
 - Unique labels
 - No overlapping or embedding
 - Relations between number of arguments; order constraints
 - If verb is of type A, no argument of type B

Constraints: examples

- No duplicate argument classes: $\sum_{i=1}^{n} a_{i,Arg0} \leq 1$
- On discontinuous arguments (C-ARG) $\forall j (1 \le j \le n), \sum_{i=1}^{j-1} a_{i,Arg0} \ge a_{j,C-Arg0}$
- On reference arguments (R-ARG)
 ∀j(1 ≤ j ≤ n), ∑_{i≠j} a_{i,Arg0} ≥ a_{j,R-Arg0}
- Many other possible constraints:
 - Unique labels
 - No overlapping or embedding
 - Relations between number of arguments; order constraints
 - If verb is of type A, no argument of type B

Constraints: examples

- No duplicate argument classes: $\sum_{i=1}^{n} a_{i,Arg0} \leq 1$
- On discontinuous arguments (C-ARG) $\forall j (1 \le j \le n), \sum_{i=1}^{j-1} a_{i,Arg0} \ge a_{j,C-Arg0}$
- On reference arguments (R-ARG)

[The deregulation]_{*Arg1*} of railroads and trucking companies [that]_{*R*-*Arg1*} began [in 1980]_{*AM*-*TMP*} enabled ... $\forall j(1 \le j \le n), \sum_{i \ne j} a_{i,Arg0} \ge a_{i,R-Arg0}$

- Many other possible constraints:
 - Unique labels
 - No overlapping or embedding
 - Relations between number of arguments; order constraints
 - If verb is of type A, no argument of type B

Constraints: examples

- No duplicate argument classes: $\sum_{i=1}^{n} a_{i,Arg0} \leq 1$
- On discontinuous arguments (C-ARG) $\forall j (1 \le j \le n), \sum_{i=1}^{j-1} a_{i,Arg0} \ge a_{j,C-Arg0}$
- On reference arguments (R-ARG) $\forall j (1 \le j \le n), \sum_{i \ne j} a_{i,Arg0} \ge a_{j,R-Arg0}$
- Many other possible constraints:
 - Unique labels
 - No overlapping or embedding
 - Relations between number of arguments; order constraints
 - If verb is of type A, no argument of type B

Constraints: examples

- No duplicate argument classes: $\sum_{i=1}^{n} a_{i,Arg0} \leq 1$
- On discontinuous arguments (C-ARG) $\forall j (1 \le j \le n), \sum_{i=1}^{j-1} a_{i,Arg0} \ge a_{j,C-Arg0}$
- On reference arguments (R-ARG) $\forall j (1 \le j \le n), \sum_{i \ne j} a_{i,Arg0} \ge a_{j,R-Arg0}$
- Many other possible constraints:
 - Unique labels
 - No overlapping or embedding
 - Relations between number of arguments; order constraints
 - If verb is of type A, no argument of type B

Constraints: examples

- No duplicate argument classes: $\sum_{i=1}^{n} a_{i,Arg0} \leq 1$
- On discontinuous arguments (C-ARG) $\forall j (1 \le j \le n), \sum_{i=1}^{j-1} a_{i,Arg0} \ge a_{j,C-Arg0}$
- On reference arguments (R-ARG) $\forall j (1 \le j \le n), \sum_{i \ne j} a_{i,Arg0} \ge a_{j,R-Arg0}$
- Many other possible constraints:
 - Unique labels
 - No overlapping or embedding
 - Relations between number of arguments; order constraints
 - If verb is of type A, no argument of type B



- Joint inference improves results $> 2.0 F_1$ points
- $\bullet\,$ Inference with many parsers improves results $\sim 2.6\ F_1$ points
- Best results at CoNLL-2005 shared task (Carreras & Màrquez, 2005)



Detour to Machine Learning Concepts (II)

What have we used from ML now?

- Inference with local classifiers under structural and problem-dependent constraints (CSP)
- Integer Linear Programming formulation
 - Efficient ILP (exact) solvers exist
 - Example: Joint learning of named entities and relations



Architecture

- Use a probabilistic local SRL model to produce multiple (*n*-best) candidate solutions for the predicate structure
- Use a feature-rich reranking model to select the best solution among them

Main goal: is to build a rich model for joint scoring, which takes into account the dependencies among the labels of argument phrases



Local Steps

- i. Parse the sentence and apply pruning (Xue & Palmer, 2004) to filter argument candidates for a given predicate p
- ii. Apply a simple local scoring model trained with log-linear classifiers (MaxEnt): P(label_i|node, p) probability distribution
- iii. Consider a simple global scoring scheme assuming independence of local assignments: $P_{LOCAL}(L|tree, p) = \prod_{node_i \in tree} P(label_i|node_i, p)$
- iv. Use dynamic programming to find the n-most probable non-overlapping complete labelings for predicate p



Reranking Step

 i. Consider a reranking model trained to select the best among the *n*-most probable complete labelings; again a log-linear model: *P_{JOINT}(L_i|tree, p)*

ii. Consider the following combination of local and joint scoring models: log(P_{SRL}(L|tree, p)) = log(P_{JOINT}(L|tree, p)) + \lambda log(P_{LOCAL}(L|tree, p))

iii. Select the complete labeling $(L_i \in \{L_1, L_2, ..., L_n\})$ that maximizes the previous formula (reranking)





Joint Model Features



Repetition features: count of arguments with a given label c(AM-TMP)=2

Complete sequence syntactic-semantic features for the core arguments:

[NP A0 hit NP A1], [NP A0 VBD NP A1] (backoff) [NP A0 hit] (left backoff) [NP ARG hit NP ARG] (no specific labels) [1 hit 1] (counts of left and right core arguments)

66

Enhancement by using multiple trees

- For top k trees from Charniak's parser, t₁, t₂,..., t_k, find corresponding best SRL assignments L₁, L₂,..., L_k and choose the tree and assignment that maximize the score (approx. joint probability of tree and assignment) score(L_i, t_i) = αlog(P(t_i)) + log(P_{SRL}(L_i|t_i))
- Final Results (2nd best at CoNLL): WSJ-23: 78.45 (F₁), 79.54 (Prec.), 77.39 (Rec.) Brown: 67.71 (F₁), 70.24 (Prec.), 65.37 (Rec.) Bug-fixed post-evaluation: 80.32 F₁ (WSJ) 68.81 F₁ (Brown)
- Improvement due to the joint model: $>2 F_1$ points



Enhancement by using multiple trees

- For top k trees from Charniak's parser, t₁, t₂,..., t_k, find corresponding best SRL assignments L₁, L₂,..., L_k and choose the tree and assignment that maximize the score (approx. joint probability of tree and assignment) score(L_i, t_i) = αlog(P(t_i)) + log(P_{SRL}(L_i|t_i))
- Final Results (2nd best at CoNLL): WSJ-23: 78.45 (F₁), 79.54 (Prec.), 77.39 (Rec.) Brown: 67.71 (F₁), 70.24 (Prec.), 65.37 (Rec.) Bug-fixed post-evaluation: 80.32 F₁ (WSJ) 68.81 F₁ (Brown)
- Improvement due to the joint model: $>2 F_1$ points



Enhancement by using multiple trees

- For top k trees from Charniak's parser, t₁, t₂,..., t_k, find corresponding best SRL assignments L₁, L₂,..., L_k and choose the tree and assignment that maximize the score (approx. joint probability of tree and assignment) score(L_i, t_i) = αlog(P(t_i)) + log(P_{SRL}(L_i|t_i))
- Final Results (2nd best at CoNLL): WSJ-23: 78.45 (F₁), 79.54 (Prec.), 77.39 (Rec.) Brown: 67.71 (F₁), 70.24 (Prec.), 65.37 (Rec.) Bug-fixed post-evaluation: 80.32 F₁ (WSJ) 68.81 F₁ (Brown)
- Improvement due to the joint model: $>2 F_1$ points



Semantic Role Labeling — A Running Example

Detour to Machine Learning Concepts (III)

What else do we need from ML?

Ranking and re-ranking algorithms (*learning to rank*)
 E.g., Ranking Perceptron



SRL as sequential tagging

- Explore the sentence regions defined by the clause boundaries.
- The top-most constituents in the regions are selected as tokens.
- Equivalent to (Xue&Palmer 04) pruning process on full parse trees





SRL as sequential tagging

- Overall results on development set

	F ₁	Prec.	Rec.
PPUPC	73.57	76.86	70.55
FP _{CHA}	75.75	78.08	73.54
Combined	76.93	78.39	75.53

- Final results on test sets
 - WSJ-23 (2416 sentences)
 - 77.97 (F₁), 79.55 (Prec.), 76.45 (Rec.)
 - Brown (426 sentences; cross-domain test)
 - 67.42 (F₁), 70.79 (Prec.), 64.35 (Rec.)



Detour to Machine Learning Concepts (IV)

More things to learn from Machine Learning?

- Sequential tagging/segmentation paradigm
 - HMMs (generative models)
 - Chained local classifiers, MEMMs, CRFs, structure perceptron


SRL Architecture

Exceptions to the standard architecture

- Parsing variations for SRL
 - ⇒ Syntactic parser trained to predict argument candidates (Yi & Palmer, 2005)
 - ⇒ Joint parsing and SRL: semantic parsing (Musillo & Merlo, 2006; Merlo & Musillo, 2008)
 - \Rightarrow SRL based on dependency parsing (Johansson & Nugues, 2007)
 - ⇒ Systems from the CoNLL-2008 and 2009 shared tasks (Surdeanu et al., 2008; Hajič et al., 2009)
 - \Rightarrow CCG parser (Gildea and Hockenmaier, 2005; Boxwell et al., 2009)
 - ⇒ HPSG parsers with handcrafted grammars (Zhang et al., 2008; 2009)



SRL Architecture

Exceptions to the standard architecture (II)

- SRL as sequential tagging (Hacioglu et al., 2004; Màrquez et al., 2005; Surdeanu et al., 2007)
- Joint treatment of all predicates in the sentence (Carreras et al., 2004; Surdeanu et al., 2008)
- SRL using Markov Logic Networks (Meza-Ruiz & Riedel, 2008; 2009)



Talk Overview

1 Motivation

Semantic Role Labeling — A Running Example

- The Statistical Approach to SRL
- Examples of SRL Systems

Feature Engineering

- Semantic Features for SRL
- An Arc-factored Model for Joint Syntactic-SRL Parsing
- Not Addressed in this Course

3 Conclusions



Features: local scoring

(Gildea & Jurafsky, 2002)

- Highly influential for the SRL work. They characterize:
 - i. The candidate argument (constituent) and its context: phrase type, head word, governing category of the constituent
 - ii. The verb predicate and its context: lemma, voice, subcategorization pattern of the verb
 - iii. The relation between the consituent and the predicate: position of the constituent with respect to the verb, category path between them.



Features: local scoring — extensions

- "Brute force" features. Applied to the constituent and possibly to parent and siblings:
 - ⇒ First and last words/POS in the constituent, bag-of-words, *n*-grams of POS, and sequence of top syntactic elements in the constituent.
- Linguistically-inspired features

⇒ Content word, named entities (Surdeanu et al., 2003), syntactic frame (Xue & Palmer, 2004), path variations, semantic compatibility between constituent head and predicate (Zapirain et al., 2007; 2009), etc.

• Significant (and cumulative) increase in performance



Features: local scoring — extensions

- "Brute force" features. Applied to the constituent and possibly to parent and siblings:
 - ⇒ First and last words/POS in the constituent, bag-of-words, *n*-grams of POS, and sequence of top syntactic elements in the constituent.
- Linguistically-inspired features
 - ⇒ Content word, named entities (Surdeanu et al., 2003), syntactic frame (Xue & Palmer, 2004), path variations, semantic compatibility between constituent head and predicate (Zapirain et al., 2007; 2009), etc.
- Significant (and cumulative) increase in performance



Features: local scoring — extensions

- "Brute force" features. Applied to the constituent and possibly to parent and siblings:
 - ⇒ First and last words/POS in the constituent, bag-of-words, *n*-grams of POS, and sequence of top syntactic elements in the constituent.
- Linguistically-inspired features
 - ⇒ Content word, named entities (Surdeanu et al., 2003), syntactic frame (Xue & Palmer, 2004), path variations, semantic compatibility between constituent head and predicate (Zapirain et al., 2007; 2009), etc.
- Significant (and cumulative) increase in performance



Features: joint scoring

- Richer features taking into account information from several arguments at a time
- Best example: when doing re-ranking one may codify patterns on the whole candidate argument structure (Hiaghighi et al., 2005; Toutanova et al., 2008)
- Good for capturing global preferences



Features: the Kernel approach

- Knowledge poor approach
- Let the kernel function to compute the similarity/differences between examples by considering all possible substructures as features
- Motivation: avoid intense knowledge engineering
- Potentially useful for rapid system development and working with under resourced languages
- Mostly variants of Collins' all-subtrees convolution kernel (Moschitti et al., 2008; Pighin & Moschitti, 2009; 2010)



Features: the Kernel approach

- Knowledge poor approach
- Let the kernel function to compute the similarity/differences between examples by considering all possible substructures as features
- Motivation: avoid intense knowledge engineering
- Potentially useful for rapid system development and working with under resourced languages
- Mostly variants of Collins' all-subtrees convolution kernel (Moschitti et al., 2008; Pighin & Moschitti, 2009; 2010)



Features: the Kernel approach

- Knowledge poor approach
- Let the kernel function to compute the similarity/differences between examples by considering all possible substructures as features
- Motivation: avoid intense knowledge engineering
- Potentially useful for rapid system development and working with under resourced languages
- Mostly variants of Collins' all-subtrees convolution kernel (Moschitti et al., 2008; Pighin & Moschitti, 2009; 2010)



Talk Overview

1 Motivation

Semantic Role Labeling — A Running Example

- The Statistical Approach to SRL
- Examples of SRL Systems
- Feature Engineering

Semantic Features for SRL

- An Arc-factored Model for Joint Syntactic-SRL Parsing
- Not Addressed in this Course

3 Conclusions



Joint work with

Eneko Agirre, Mihai Surdeanu and Beñat Zapirain

(Zapirain et al. 2010) — ACL (Zapirain et al. 2011) — NAACL (Zapirain et al. 2013) — Computational Linguistics 39(3)

Results from CoNLL-2005 shared task

Results on WSJ and Brown Tests





Results from CoNLL-2005 shared task

Reasons for the low generalization ability

- The training corpus is not representative and large enough (and it will never be)
- Taggers and syntactic parsers also experience a significant drop in performance
- The main loss in performance takes place in role classification, not identification — semantic explanation (Pradhan et al., 2008)



Semantic Features for SRL

Motivation

- Most current systems capture semantics through lexicalized features on the predicate and the head word of the argument to be classified
- But lexical features are sparse and generalize badly
 [JFK]_{Patient} was_assassinated [in Dallas]_{LOC}
 [JFK]_{Patient} was_assassinated [in November]_{TMP}
- [in Texas]???, [in autumn]???



Semantic Features for SRL

Motivation

- Most current systems capture semantics through lexicalized features on the predicate and the head word of the argument to be classified
- But lexical features are sparse and generalize badly
 [JFK]_{Patient} was_assassinated [in Dallas]_{LOC}
 [JFK]_{Patient} was_assassinated [in November]_{TMP}
- [in Texas]???, [in autumn]???



Semantic Features for SRL

Motivation

- Most current systems capture semantics through lexicalized features on the predicate and the head word of the argument to be classified
- But lexical features are sparse and generalize badly
 [JFK]_{Patient} was_assassinated [in Dallas]_{LOC}
 [JFK]_{Patient} was_assassinated [in November]_{TMP}
- [in Texas]???, [in autumn]???



Semantic Role Labeling — A Running Example

Semantic Features for SRL

Motivation

Selectional Preferences and distributional similarity techniques should help us to classify arguments with low-frequency or unknown head words

 $[Dallas \approx Texas]_{Location}, [November \approx autumn]_{Temporal}$



Previous Work

Selectional Preferences

- Modeling semantic preferences that predicates impose on their arguments
- Long tradition of automatic acquisition of selectional preferences (SPs) from corpora. WordNet-based and distributional models of SPs (Resnik, 1993; Pantel and Lin, 2000; Brockmann and Lapata, 2003) (Erk 2007; Erk et al., 2011; etc.)
 - ⇒ e.g., estimate plausibility of triples: (verb, argument, head-word)
 - \Rightarrow useful for syntactic-semantic disambiguation



Previous Work

SPs applied to Semantic Role Labeling

- (Gildea and Jurafsky, 2002) FrameNet
 - \Rightarrow First researchers to apply selectional preferences to SRL
 - ⇒ Distributional clustering and WordNet-based techniques to generalize argument heads
 - ⇒ Slight improvement in role classification (NP arguments)
- Zapirain et al. (2010; 2013) PropBank
 - ⇒ Show that selectional preferences can improve semantic role classification in a state-of-the-art SRL system



Two types of selectional preferences (SP)

i. **verb**-*role*: list of heads of NP arguments of the predicate **verb** that are labeled with the role *role*

write-Arg0: Angrist anyone baker ball bank Barlow Bates ... write-Arg1: abstract act analysis article asset bill book ... write-Arg2: bank commander hundred jaguar Kemp member ... write-AM-LOC: paper space ...

. . .

ii. **prep**-*role*: list of nominal heads of PP arguments with preposition **prep** that are labeled with the role *role*

from-Arg2: academy account acquisition activity ad ...
from-Arg3: activity advertising agenda airport ...
from-Arg4: europe Golenbock system Vizcaya west
from-AM-TMP: april august beginning bell day dec. half ...
from-AM-LOC: agency area asia body bureau orlando ...

SP models: $SP_{sim}(p, r, w)$ compatibility score

- Discriminative approach: given a new argument of a predicate *p*, we compare its head (*w*) to the selectional preference of each possible role label *r*, i.e., we want to find the role with the selectional preference that fits the head best
- We compute the compatibility scores using two different methods
 - \Rightarrow WordNet based —using (Resnik, 1993)
 - \Rightarrow Based on distributional similarity —a la Erk (2007)



SP models: $SP_{sim}(p, r, w)$ compatibility score

- Discriminative approach: given a new argument of a predicate *p*, we compare its head (*w*) to the selectional preference of each possible role label *r*, i.e., we want to find the role with the selectional preference that fits the head best
- We compute the compatibility scores using two different methods
 - \Rightarrow WordNet based —using (Resnik, 1993)
 - \Rightarrow Based on distributional similarity —a la Erk (2007)



(Zapirain et al., 2013)

WordNet SP models

• Resnik formula (1993) is used to precalculate a weighted list of relevant synsets for the lists of words contained in the SPs

SP write-Arg0: Angrist anyone baker ball bank Barlow Bates ...

n#0002086 5.875 life form organism being living thing "any living entity" n#00001740 5.737 entity something "anything having existence (living or nonliving)" n#00009457 4.782 object physical object "a physical (tangible and visible) entity;" n#00004123 4.351 person individual someone somebody mortal human soul "a human being;" ...

SP write-Arg1: abstract act analysis article asset bill book ...

n#00019671 7.956 communication "something that is communicated between people or groups" n#04949838 4.257 message content subject matter substance "what a communication that ..." n#00018916 3.848 relation "an abstraction belonging to or characteristic of two entities" n#00013018 3.574 abstraction "a concept formed by extracting common features from examples"



WordNet SP models

- At test time, for a new argument of the predicate **write** with head word book:
 - ⇒ consider S = {<book>} ∪ "all its hypernyms in WordNet" (for all senses of book)
 - ⇒ SP_{Res}(write, Arg1, book) returns the sum of the weights of the sysnsets in S matching the synsets in the list corresponding to the SP write-Arg1



Distributional SP models: based on Erk's (2007) setting

JFK was assassinated [in Texas] ???

SP in-*TMP*: November, century, month SP in-*LOC*: Dallas, railway, city

 $SP_{sim}(p, r, w) = \sum_{w_i \in Seen(p, r)} sim(w, w_i) \cdot weight(p, r, w_i)$

SP(in, TMP, Texas) = sim(Texas, November) · weight(in, TMP, November) +
sim(Texas, century) · weight(in, TMP, century) +
sim(Texas, month) · weight(in, TMP, month)



Distributional SP models: based on Erk's (2007) setting

JFK was assassinated [in Texas] ???

SP in-*TMP*: November, century, month SP in-*LOC*: Dallas, railway, city

$$SP_{sim}(p, r, w) = \sum_{w_i \in Seen(p, r)} sim(w, w_i) \cdot weight(p, r, w_i)$$

SP(in, TMP, Texas) = sim(Texas, November) · weight(in, TMP, November) + sim(Texas, century) · weight(in, TMP, century) + sim(Texas, month) · weight(in, TMP, month)



Distributional SP models: based on Erk's (2007) setting

JFK was assassinated [in Texas] ???

SP in-*TMP*: November, century, month SP in-*LOC*: Dallas, railway, city

$$SP_{sim}(p, r, w) = \sum_{w_i \in Seen(p, r)} sim(w, w_i) \cdot weight(p, r, w_i)$$

SP(in, TMP, Texas) = sim(Texas, November) · freq(in, TMP, November) + sim(Texas, century) · freq(in, TMP, century) + sim(Texas, month) · freq(in, TMP, month)



(Zapirain et al., 2013)

Distributional SP models: based on Erk's (2007) setting

JFK was assassinated [in Texas]???

SP in-*TMP*: November, century, month SP in-*LOC*: Dallas, railway, city

$$SP_{sim}(p, r, w) = \sum_{w_i \in Seen(p, r)} sim(w, w_i) \cdot weight(p, r, w_i)$$

SP(in, LOC, Texas) = sim(Texas, Dallas) · freq(in, LOC, Dallas) + sim(Texas, railway) · freq(in, LOC, railway) + sim(Texas, city) · freq(in, LOC, city)

SP(in,LOC,Texas) > SP(in,TMP,Texas)

anar compoong MESEBER I I STILLE desht of Qate Jandatan Jai-baada, gada

(Zapirain et al., 2013)

Distributional SP models: based on Erk's (2007) setting

JFK was assassinated [in Texas]???

SP in-*TMP*: November, century, month SP in-*LOC*: Dallas, railway, city

$$SP_{sim}(p, r, w) = \sum_{w_i \in Seen(p, r)} sim(w, w_i) \cdot weight(p, r, w_i)$$

SP(in, LOC, Texas) = sim(Texas, Dallas) · freq(in, LOC, Dallas) + sim(Texas, railway) · freq(in, LOC, railway) + sim(Texas, city) · freq(in, LOC, city)

SP(in,LOC,Texas) > SP(in,TMP,Texas)

Distributional SP models: various instantiations for sim

- Using Padó and Lapata's software (2007) for computing distributional similarity measures
 - \Rightarrow Run on the British National Corpus
 - \Rightarrow Optimal parameterization as described in the paper
 - ⇒ Jaccard, cosine and Lin's similarity measures: sim_{Jac}, sim_{cos} and sim_{Lin}
- Using the already available Lin's thesaurus (Lin, 1998)
 - \Rightarrow Direct and second order similarity: sim_{Lin}^{th} , sim_{Jac}^{th2} and sim_{cos}^{th2}
 - \Rightarrow Average of both directions similarity



Setting: Assign role labels to argument head words based solely on SP scores

- ⇒ For each head word (w), select the role (r) of the predicate or preposition (p) which fits best the head word: R_{sim}(p, w) = arg max_{r∈Roles}(p) SP_{sim}(p, r, w)
- \Rightarrow SPs based on (p, r, w) triples from CoNLL-2005 data
- ⇒ In-domain (WSJ) and out-of-domain (Brown) test sets CoNLL-2005
- ⇒ Lexical baseline model: for a test pair (p, w), assign the role under which the head (w) occurred most often in the training data given the predicate (p)



Setting: Assign role labels to argument head words based solely on SP scores

- ⇒ For each head word (w), select the role (r) of the predicate or preposition (p) which fits best the head word: $R_{sim}(p, w) = \arg \max_{r \in Roles(p)} SP_{sim}(p, r, w)$
- \Rightarrow SPs based on (p, r, w) triples from CoNLL-2005 data
- ⇒ In-domain (WSJ) and out-of-domain (Brown) test sets CoNLL-2005
- ⇒ Lexical baseline model: for a test pair (p, w), assign the role under which the head (w) occurred most often in the training data given the predicate (p)



Setting: Assign role labels to argument head words based solely on SP scores

- ⇒ For each head word (w), select the role (r) of the predicate or preposition (p) which fits best the head word: $R_{sim}(p, w) = \arg \max_{r \in Roles(p)} SP_{sim}(p, r, w)$
- \Rightarrow SPs based on (p, r, w) triples from CoNLL-2005 data
- \Rightarrow In-domain (WSJ) and out-of-domain (Brown) test sets CoNLL-2005
- ⇒ Lexical baseline model: for a test pair (p, w), assign the role under which the head (w) occurred most often in the training data given the predicate (p)



	WSJ-test			Brown		
	prec.	rec.	F_1	prec.	rec.	F_1
lexical	82.98	43.77	57.31	68.47	13.60	22.69
SP_{Res}	63.47	53.24	57.91	55.12	44.15	49.03
$SP_{sim_{Jac}}$	61.83	61.40	61.61	55.42	53.45	54.42
$SP_{sim_{cos}}$	64.67	64.22	64.44	56.56	54.54	55.53
$SP_{sim_{las}^{th2}}$	70.82	70.33	70.57	62.37	60.15	61.24
SP _{simth2}	70.28		70.04	62.36	60.14	61.23

\Rightarrow Lexical features have a high precision but very low recall

- \Rightarrow SPs are able to effectively generalize lexical features
- \Rightarrow SPs based on distributional similarity are better
- \Rightarrow Second-order similarity variants (Lin) attain the best results

anan composing Research Institute Mindre of Quar Mondation Juli Acade, pipela
Evaluation of SPs in isolation

		WSJ-test			Brown	
	prec.	rec.	F_1	prec.	rec.	F_1
lexical	82.98	43.77	57.31	68.47	13.60	22.69
SP_{Res}	63.47	53.24	57.91	55.12	44.15	49.03
$SP_{sim_{Jac}}$	61.83	61.40	61.61	55.42	53.45	54.42
$SP_{sim_{cos}}$	64.67	64.22	64.44	56.56	54.54	55.53
$SP_{sim_{lac}^{th2}}$	70.82	70.33	70.57	62.37	60.15	61.24
$SP_{sim_{cos}^{th2}}$	70.28	69.80	70.04	62.36	60.14	61.23

- \Rightarrow Lexical features have a high precision but very low recall
- \Rightarrow SPs are able to effectively generalize lexical features
- \Rightarrow SPs based on distributional similarity are better
- \Rightarrow Second-order similarity variants (Lin) attain the best results

- SwiRL system for SRL (Surdeanu et al., 2007)
 - ⇒ System from CoNLL-2005 shared task (PropBank)
 - \Rightarrow Standard architecture (ML based on AdaBoost and SVMs)
 - ⇒ Best results from single (non-combined) systems at CoNLL-2005
- Simple approach: extending *SwiRL* features with SP predictions
 - \Rightarrow We train several extended *SwiRL-SP*_i models, one per selectional preferences model *SP*_i
 - ⇒ For each example (p, w) of *SwiRL-SP_i*, we add a single new feature whose value is the predicted role label $R_i(p, w)$



Results

	WSJ-test			Brown		
	Core	Adj	All	Core	Adj	All
SwiRL	93.25	81.31	90.83	84.42	57.76	79.52
SwiRL+SP _{Res}	93.17	81.08	90.76	84.52	59.24	79.86
SwiRL+SP _{simJac}	93.37	80.30	90.86	84.43	59.54	79.83
SwiRL+SP _{simcos}	93.33	80.92	90.87	85.14	60.16	80.50
SwiRL+SP _{sim^{th2}}	93.03	82.75	90.95	85.62	59.63	80.75
$SwiRL+SP_{sim_{cos}^{th2}}$	93.78	80.56	91.23	84.95	61.01	80.48

- $\Rightarrow\,$ Slight improvements, especially noticeable on Brown corpus
- \Rightarrow Weak signal of a single feature?



- Simple combinations of the individual *SwiRL*+*SP_i* classifiers worked quite well (majority voting)
- We also trained a meta-classifier to combine the *SwiRL*+*SP_i* classifiers and the stand-alone *SP_i* models:
 - ⇒ Binary classification approach: "is a proposed role correct or not?"
 - \Rightarrow Features are based on the predictions of base SP_i and $SwiRL+SP_i$ models
 - $\Rightarrow\,$ Trained with a SVM with a quadratic polynomial kernel



Results (II)

	WSJ-test			Brown		
	Core	Adj	All	Core	Adj	All
SwiRL	93.25	81.31	90.83	84.42	57.76	79.52
$+SP_{sim_{cos}^{th2}}$	93.78	80.56	91.23	84.95	61.01	80.48
Meta	94.37	83.40	92.12	86.20	63.40	81.91

• Statistically significant improvements (99%) for both core and adjunct arguments, both in domain and out of domain



Results (II)

	WSJ-test			Brown		
	Core	Adj	All	Core	Adj	All
SwiRL	93.25	81.31	90.83	84.42	57.76	79.52
$+SP_{sim_{cos}^{th2}}$	93.78	80.56	91.23	84.95	61.01	80.48
Meta	94.37	83.40	92.12	86.20	63.40	81.91

• Statistically significant improvements (99%) for both core and adjunct arguments, both in domain and out of domain



Output analysis

- Manual inspection of 50 cases in which the meta classifier corrects SwiRL:
 - $\Rightarrow\,$ Usually cases with low frequency verbs or argument heads
 - ⇒ In ~58% of the cases, syntax does not disambiguate, seems to suggest a wrong role label or it is confusing SwiRL because it is incorrect. However, most of the SP predictions are correct.
 - ⇒ ~30% of the cases: unclear source of the SwiRL error but still several SP models suggest the correct role
 - \Rightarrow \sim 12% of the cases: chance effect



Output analysis

- Manual inspection of 50 cases in which the meta classifier corrects SwiRL:
 - \Rightarrow Usually cases with low frequency verbs or argument heads
 - ⇒ In ~58% of the cases, syntax does not disambiguate, seems to suggest a wrong role label or it is confusing SwiRL because it is incorrect. However, most of the SP predictions are correct.
 - $\Rightarrow~{\sim}30\%$ of the cases: unclear source of the SwiRL error but still several SP models suggest the correct role
 - \Rightarrow ~12% of the cases: chance effect



С	utput	ana	lysis:	examp	le 1
			2		

		Several	JJ	(S1(S(NP*
		traders	NNS	*)
		could	MD	(VP*
		be	VB	(VP*
		seen	VBN	(VP*
		shaking	VBG	(S(VP*
		their	PRP\$	(NP*
		heads	NNS	*)))
		when	WRB	(SBAR(WHADVP*)
A1	A0	the	DT	(S(NP*
A1	A0	news	NN	*)
	(P)	flashed	VBD	(VP*))))))
				*))

|--|

		Italian	NNP	(S1(S(NP*
		President	NNP	*
		Francesco	NNP	*
		Cossiga	NNP	*)
	(P)	promised	VBD	(VP* [′]
12	A1	а	DT	(NP(NP*
12	A1	quick	JJ	*
12	A1	investigation	NN	*)
12	A1	into	IN	(PP*́
12	A1	whether	IN	(SBAR*
12	A1	Olivetti	NNP	(S(NP*)
12	A1	broke	VBD	`(VP*́
12	A1	Cocom	NNP	(NP*
12	A1	rules	NNS	` *)))))))
				*))

Output analysis: example 3

	(P)
A3	TMP

Annual	JJ	(S(NP*
payments	NNS	*)
will	MD	(VP*
more	RBR	(VP(ADVP*
than	IN	*)
double	VB	*
from	IN	(PP*
а	DT	(NP*
year	NN	`*
ago	RB	*))
to	ТО	(PP*
about	RB	(NP(QP*
\$240	CD	*
million	CD	*)))
		,,,

n montute Geologia de la de

Output analysis: example 4

		Procter	NNP	(S1(S(NP*
		&	CC	*
		Gamble	NNP	*
		Co.	NNP	*)
		plans	VBZ	(VP*
		to	то	(S(VP*
		begin	VB	(VP*
	(P)	testing	VBG	(S(VP*
		next	JJ	(NP*
		month	NN	*)))
A1	A0	а	DT	(NP(NP*
A1	A0	superco.	JJ	*
A1	A0	detergent	NN	*)
A1	A0	that	WDT	(SBAR(WHNP*)
A1	A0	washload	NN	(NP*))))))))))))))))))))))))))))))))))))
				*))

Final Remarks...

...on Semantic Features and Generalizations for SRL

- Word Embeddings (and learning with deep NNs)
 - Turian et al., ACL 2010
 - Collobert et al., JMLR 2011 (SENNA)
 - Foland and Martin, NAACL 2015
- Low-rank decomposition of high-order tensor models
 - Lei et al., NAACL 2015



Talk Overview

1 Motivation

Semantic Role Labeling — A Running Example

- The Statistical Approach to SRL
- Examples of SRL Systems
- Feature Engineering
- Semantic Features for SRL
- An Arc-factored Model for Joint Syntactic-SRL Parsing
- Not Addressed in this Course

3 Conclusions



Joint work with

Xavier Lluís and Xavier Carreras

(Lluís et al. 2013) — TACL (to be presented at ACL)

CoNLL-2008/2009 shared task

Joint parsing of syntactic and semantic dependencies





An Arc-factored Model for Joint Syntactic-SRL Parsing

A Simplified Example



• Predicate-argument structures are naturally represented with dependencies



A Simplified Example



- Semantic roles are strongly related to syntactic structure
- Typical systems find semantic roles in a pipeline
 - \Rightarrow First obtain the syntactic tree
 - $\Rightarrow\,$ Second obtain the semantic roles, using the syntactic tree
- Pipeline systems can not correct syntax based on semantic roles



A Simplified Example



- We model the two structures jointly
 - ⇒ To capture interactions between syntactic and semantic dependencies
- Challenge:
 - $\Rightarrow\,$ Some semantic dependencies are associated with a segment of syntactic dependencies
 - \Rightarrow Hard to factorize the two structures jointly



Decomposing Syntactic and Semantic Trees



Syntactic subproblem

$$\begin{array}{rcl} \textit{syn}(\mathbf{x}) & = & \operatorname*{argmax}_{\mathbf{y}} \; \text{score}_\textit{syn}(\mathbf{x},\mathbf{y}) \\ & & \underset{\mathbf{y}}{\text{subject to}} \; \; \mathbf{cTree:} \; \mathbf{y} \; \text{is a projective tree} \end{array}$$

- Solved by a standard dependency parsing algorithm
- score_syn(\mathbf{x}, \mathbf{y}) is arc-factored: 1st and 2nd order models
- Graph-based parsing algorithms, reimplementing (McDonald, 2005; Carreras et al., 2007)
- Trained with (linear) average structure perceptron using state-of-the-art features



Semantic Subproblem

 $\begin{array}{lll} \mathit{srl}(\mathbf{x}) & = & \underset{\mathbf{z},\pi}{\operatorname{argmax}} \operatorname{score_srl}(\mathbf{x},\mathbf{z},\pi) \\ & \text{subject to} & \mathsf{cRole:} \text{ no repeated roles} \\ & \mathsf{cArg:} \text{ at most one role per token} \\ & \mathsf{cPath:} \ \pi \ \mathsf{codifies \ paths \ consistent \ with \ \mathbf{z}} \end{array}$

- In a predicate:
 - \Rightarrow A token appears at most once as argument
 - \Rightarrow A semantic role appears at most once
- score_srl($\mathbf{x}, \mathbf{z}, \pi$) is factorized at the level of $\langle \mathbf{x}, p, a, r, \pi^{p,a,r} \rangle$
- local score_srl($\mathbf{x}, p, a, r, \pi^{p,a,r}$) provided by linear classifiers
- We frame the argmax inference as a linear assignment problem



SRL as Assignment



- The Hungarian algorithm solves it in $\mathcal{O}(n^3)$
- $w_{i,j}$ are the previous local predictions score_srl($\mathbf{x}, p, a, r, \pi^{p,a,r}$)
- In practice, the list of most likely paths from *p* to *a* is pre-computed using syntactic models
- Learning is performed with structure perceptron, with feedback applied after solving the assignment problem



An Arc-factored Model for Joint Syntactic-SRL Parsing

Joint Syntactic-Semantic Inference

$$\begin{array}{ll} \langle y^*,z^*,\pi^*\rangle & = & \underset{y,z,\pi}{\operatorname{argmax}} sc_syn(x,y) + sc_srl(x,z,\pi) \\ & \text{subject to} & c\mathsf{Tree, cRole, cArg, cPath} \\ & & c\mathsf{Subtree:} \ y \ \text{is consistent with } \pi \end{array}$$



Joint Syntactic-Semantic Inference

$$\begin{array}{ll} \langle y^{*},z^{*},\pi^{*}\rangle & = & \underset{y,z,\pi}{\operatorname{argmax}} \operatorname{sc_syn}(x,y) + \operatorname{sc_srl}(x,z,\pi) \\ & \text{subject to} & \operatorname{cTree, cRole, cArg, cPath} \\ & & \operatorname{cSubtree:} y \text{ is consistent with } \pi \end{array}$$

cSubtree constraints can be easily expressed as:

$$orall d \in \mathbf{y}$$
 , $c \cdot \mathbf{y}_d \geqslant \sum_{p, a, r \in \mathbf{z}} oldsymbol{\pi}_d^{p, a, r}$

or, equivalently, as equality constraints

$$orall d \in \mathbf{y}$$
 , $c \cdot \mathbf{y}_d - \sum_{p,a,r \in \mathbf{z}} \pi_d^{p,a,r} - \xi_d = 0$



Joint Syntactic-Semantic Inference

- We employed Dual Decomposition to solve the joint inference (Rush and Collins, 2011) (Sontag et al 2010)
- Lagrangian relaxation-based method that iteratively solves decomposed sub-problems with agreement constraints:
 - \Rightarrow Subtree constraints are relaxed by introducing Lagrange multipliers for every dependency λ_d
 - $\Rightarrow~$ Subproblems now depend on the λ penalty variables but can be efficiently solved
 - ⇒ Syntax: standard dependency parsing inference
 - ⇒ Semantic: linear assignment
- Guaranteed optimal solution when it converges
- In experiments, convergence in > 99.5% of sentences



We ran experiments on the CoNLL-2009 datasets with the following configurations:

Pipeline best *syn* then best *srl* enforcing cArg +Assignment enforces cRole, cArg over best *syn* Forest works with a forest of *syn* trees DD applies dual-decomposition



	syn		sem	
system	acc	prec	rec	F_1
Pipeline-1				
+Assignment-1				
Forest-1				
DD-1				

Results on WSJ development set



	syn		sem	
system	асс	prec	rec	F ₁
Pipeline-1	85.32	86.23	67.67	75.83
+Assignment-1	85.32	84.08	71.82	77.47
Forest-1				
DD-1				

+Assignment improves over Pipeline



	syn		sem	
system	асс	prec	rec	F_1
Pipeline-1	85.32	86.23	67.67	75.83
+Assignment-1	85.32	84.08	71.82	77.47
Forest-1	85.32	80.67	73.60	76.97
DD-1				

Forests shows higher recall



	syn		sem	
system	асс	prec	rec	F_1
Pipeline-1	85.32	86.23	67.67	75.83
+Assignment-1	85.32	84.08	71.82	77.47
Forest-1	85.32	80.67	73.60	76.97
DD-1	85.48	83.99	72.69	77.94

DD-1 achieves better sem F_1



	syn		sem		
system	асс	prec	rec	F_1	
Pipeline-1	85.32	86.23	67.67	75.83	
+Assignment-1	85.32	84.08	71.82	77.47	
Forest-1	85.32	80.67	73.60	76.97	
DD-1	85.48	83.99	72.69	77.94	
Pipeline-2	87.77	87.07	68.65	76.77	
+Assignment-2	87.77	85.21	73.41	78.87	
Forest-2	87.77	80.67	73.60	76.97	
DD-2	87.84	85.20	73.23	78.79	

Second-order paths are quite accurate



	syn		sem				
WSJ	асс	prec	rec	F_1	PP		
Lluís09 Merlo09	87.48 88.79	73.87 81.00	67.40 76.45	70.49 78.66	39.68 54.80		
DD-2	89.21	86.01	74.84	80.04	55.73		

Results in WSJ corpus (in-domain) test set



	syn		sem			
WSJ	асс	prec	rec	F_1	PP	
Lluís09 Merlo09	87.48 88.79	73.87 81.00	67.40 76.45	70.49 78.66	39.68 54.80	
DD-2	89.21	86.01	74.84	80.04	55.73	

Better results than Merlo09



	syn		sem				
Brown	асс	prec	rec	F_1	PP		
Lluís09 Merlo09	80.92 80.84	62.29 68.97	59.22 63.06	60.71 65.89	29.79 38.92		
DD-2	82.61	74.12	61.59	67.83	38.92		

Results in Brown corpus (out-of-domain) test set



Talk Overview

1 Motivation

Semantic Role Labeling — A Running Example

- The Statistical Approach to SRL
- Examples of SRL Systems
- Feature Engineering
- Semantic Features for SRL
- An Arc-factored Model for Joint Syntactic-SRL Parsing
- Not Addressed in this Course

3 Conclusions


Other Important Topics

Learning with latent variables/structures

- Henderson et al., Computational Linguistics 39(4), 2013
- Onsupervised models for SRL
 - Titov and Khoddam, NAACL 2015
- Learning with weak/distant supervision
- Deep NN Learning for SRL
 - Collobert et al., JMLR 2011 (SENNA)
 - Foland and Martin, NAACL 2015



Talk Overview



Semantic Role Labeling — A Running Example





- NLP technology is very important for a number of current applications:
 - $\Rightarrow\,$ MT, personal assistants, information search and analysis, market study, trends, opinions, etc.
- NLP current approaches are empirical
 ⇒ based on data, statistics, and machine learnin
- ML is at many stages of NLP state-of-the-art solutions
 and it is here to stay...



- NLP technology is very important for a number of current applications:
 - $\Rightarrow\,$ MT, personal assistants, information search and analysis, market study, trends, opinions, etc.
- NLP current approaches are empirical
 - \Rightarrow based on data, statistics, and machine learning
- ML is at many stages of NLP state-of-the-art solutions
 and it is here to stay...



- NLP technology is very important for a number of current applications:
 - ⇒ MT, personal assistants, information search and analysis, market study, trends, opinions, etc.
- NLP current approaches are empirical
 - \Rightarrow based on data, statistics, and machine learning
- ML is at many stages of NLP state-of-the-art solutions
 and it is here to stay...



 If you want to work for Google, Facebook, Yahoo, Twitter, MSR, IBM Watson...

But also Bloomberg, Goldman Sachs, Machine Zone, etc.

- You better learn about:
 - ⇒ Computational Linguistics, Statistics, Machine Learning, Text mining and analysis, etc.
- ...and you conduct a PhD first (nice opportunities at the moment)



- If you want to work for Google, Facebook, Yahoo, Twitter, MSR, IBM Watson...
 But also Bloomberg, Goldman Sachs, Machine Zone, etc.
- You better learn about:
 - ⇒ Computational Linguistics, Statistics, Machine Learning, Text mining and analysis, etc.
- ...and you conduct a PhD first (nice opportunities at the moment)



- If you want to work for Google, Facebook, Yahoo, Twitter, MSR, IBM Watson...
 But also Bloomberg, Goldman Sachs, Machine Zone, etc.
- You better learn about:
 - ⇒ Computational Linguistics, Statistics, Machine Learning, Text mining and analysis, etc.
- ...and you conduct a PhD first (nice opportunities at the moment)



- If you want to work for Google, Facebook, Yahoo, Twitter, MSR, IBM Watson...
 But also Bloomberg, Goldman Sachs, Machine Zone, etc.
- You better learn about:
 - ⇒ Computational Linguistics, Statistics, Machine Learning, Text mining and analysis, etc.
- ...and you conduct a PhD first (nice opportunities at the moment)



- I opted for taking a complex enough NLP task (SRL) as an excuse to cover all possible NLP-ML topics
- We have overviewed many important concepts and methods of Machine Learning for NLP (especially supervised)
- But there are MANY MORE that we left untouched
 ⇒ some of them currently very TRENDY!



- I opted for taking a complex enough NLP task (SRL) as an excuse to cover all possible NLP-ML topics
- We have overviewed many important concepts and methods of Machine Learning for NLP (especially supervised)
- But there are MANY MORE that we left untouched
 ⇒ some of them currently very TRENDY!



- I opted for taking a complex enough NLP task (SRL) as an excuse to cover all possible NLP-ML topics
- We have overviewed many important concepts and methods of Machine Learning for NLP (especially supervised)
- But there are MANY MORE that we left untouched
 ⇒ some of them currently very TRENDY!



on Semantic Role Labeling

- SRL is an important problem in NLP, strongly related to applications requiring some degree of semantic interpretation
- It is an active topic of research, which has generated an important body of work in the last 10 years
 ⇒ techniques, resources, applications

Some news are good but...

- ⇒ SRL still has to resolve important problems before we see a spread usage in real open-domain applications
- ⇒ A jump is needed from the laboratory conditions to the real world.



on Semantic Role Labeling

- SRL is an important problem in NLP, strongly related to applications requiring some degree of semantic interpretation
- It is an active topic of research, which has generated an important body of work in the last 10 years
 ⇒ techniques, resources, applications

Some news are good but...

- ⇒ SRL still has to resolve important problems before we see a spread usage in real open-domain applications
- \Rightarrow A jump is needed from the laboratory conditions to the real world.





I hope you enjoyed this part of the course!!!



Machine Learning applied to Natural Language Processing

Lluís Màrquez

ALT Research Group Qatar Computing Research Institute

Advanced Methods for Corpus Analysis LCT – European Masters Program in Language and Communication Technologies Donostia, June 22-24, 2015