

Perplexity of n-gram and dependency language models

Martin Popel, David Mareček
ÚFAL, Charles University in Prague



TSD, 13th International Conference on Text, Speech and Dialogue
September 8, 2010, Brno

Outline

- Language Models (LM)
 - basics
 - design decisions
- Post-ngram LM
- Dependency LM
- Evaluation
- Conclusion & future plans

Language Models – basics

$P(s) = ?$

$P(\text{The dog barked again})$

Language Models – basics

$P(s) = ?$

$P(\text{The dog barked again}) > P(\text{The dock barked again})$

Language Models – basics

$$P(s) = P(w_1, w_2, \dots, w_m)$$

$$P(\text{The dog barked again}) =$$

$$P(w_1 = \text{The}, w_2 = \text{dog}, w_3 = \text{barked}, w_4 = \text{again})$$

Language Models – basics

$$P(s) = P(w_1, w_2, \dots, w_m) = P(w_1) P(w_2|w_1) \dots P(w_m|w_1, \dots, w_{m-1})$$



Chain rule

$$P(\text{The dog barked again}) =$$

$$P(w_1 = \text{The}) \cdot$$

$$P(w_2 = \text{dog} \mid w_1 = \text{The}) \cdot$$

$$P(w_3 = \text{barked} \mid w_1 = \text{The}, w_2 = \text{dog}) \cdot$$

$$P(w_4 = \text{again} \mid w_1 = \text{The}, w_2 = \text{dog}, w_3 = \text{barked})$$

Language Models – basics

$$P(s) = P(w_1, w_2, \dots, w_m) = P(w_1) P(w_2|w_1) \dots P(w_m|w_1, \dots, w_{m-1})$$

Changed notation

$$P(\text{The dog barked again}) =$$

$$P(w_i = \text{The} \mid i=1) \cdot$$

$$P(w_i = \text{dog} \mid i=2, w_{i-1} = \text{The}) \cdot$$

$$P(w_i = \text{barked} \mid i=3, w_{i-2} = \text{The}, w_{i-1} = \text{dog}) \cdot$$

$$P(w_i = \text{again} \mid i=4, w_{i-3} = \text{The}, w_{i-2} = \text{dog}, w_{i-1} = \text{barked})$$

Language Models – basics

$$P(s) = P(w_1, w_2, \dots, w_m) = P(w_1) P(w_2|w_1) \dots P(w_m|w_1, \dots, w_{m-1})$$

Artificial start-of-sentence token

$$P(\text{The dog barked again}) =$$

$$P(w_i = \text{The} \mid i=1, w_{i-1} = \text{NONE}) \cdot$$

$$P(w_i = \text{dog} \mid i=2, w_{i-2} = \text{NONE}, w_{i-1} = \text{The}) \cdot$$

$$P(w_i = \text{barked} \mid i=3, w_{i-3} = \text{NONE}, w_{i-2} = \text{The}, w_{i-1} = \text{dog}) \cdot$$

$$P(w_i = \text{again} \mid i=4, w_{i-4} = \text{NONE}, w_{i-3} = \text{The}, w_{i-2} = \text{dog}, w_{i-1} = \text{barked})$$

Language Models – basics

$$P(s) = P(w_1, w_2, \dots, w_m) = P(w_1) P(w_2|w_1) \dots P(w_m|w_1, \dots, w_{m-1})$$

Position backoff

$$P(\text{The dog barked again}) \approx$$

$$P(w_i = \text{The} \mid w_{i-1} = \text{NONE}) \cdot$$

$$P(w_i = \text{dog} \mid w_{i-2} = \text{NONE}, w_{i-1} = \text{The}) \cdot$$

$$P(w_i = \text{barked} \mid w_{i-3} = \text{NONE}, w_{i-2} = \text{The}, w_{i-1} = \text{dog}) \cdot$$

$$P(w_i = \text{again} \mid w_{i-4} = \text{NONE}, w_{i-3} = \text{The}, w_{i-2} = \text{dog}, w_{i-1} = \text{barked})$$

Language Models – basics

$$P(s) = P(w_1, w_2, \dots, w_m) \approx \prod_{i=1..m} P(w_i | w_{i-1})$$

History backoff (bigram LM)

$$P(\text{The dog barked again}) \approx$$

$$P(w_i = \text{The} \mid w_{i-1} = \text{NONE}) \cdot$$

$$P(w_i = \text{dog} \mid w_{i-1} = \text{The}) \cdot$$

$$P(w_i = \text{barked} \mid w_{i-1} = \text{dog}) \cdot$$

$$P(w_i = \text{again} \mid w_{i-1} = \text{barked})$$

Language Models – basics

$$P(s) = P(w_1, w_2, \dots, w_m) \approx \prod_{i=1..m} P(w_i | w_{i-2}, w_{i-1})$$

History backoff (trigram LM)

$$P(\text{The dog barked again}) \approx$$

$$P(w_i = \text{The} \mid w_{i-2} = \text{NONE}, w_{i-1} = \text{NONE}) \cdot$$

$$P(w_i = \text{dog} \mid w_{i-2} = \text{NONE}, w_{i-1} = \text{The}) \cdot$$

$$P(w_i = \text{barked} \mid w_{i-2} = \text{The}, w_{i-1} = \text{dog}) \cdot$$

$$P(w_i = \text{again} \mid w_{i-2} = \text{dog}, w_{i-1} = \text{barked})$$

Language Models – basics

$$P(s) = P(w_1, w_2, \dots, w_m) \approx \prod_{i=1..m} P(w_i | w_{i-2}, w_{i-1})$$

In general: $\prod_{i=1..m} P(w_i | \mathbf{h}_i)$

\mathbf{h}_i ... context (history) of word w_i

Language Models – design decisions

- 1) How to factorize $P(w_1, w_2, \dots, w_m)$ into $\prod_{i=1..m} P(w_i | \mathbf{h}_i)$,
i.e. what word-positions will be used as the context \mathbf{h}_i ?
- 2) What additional context information will be used
(apart from word forms),
e.g. stems, lemmata, POS tags, word classes,... ?
- 3) How to estimate $P(w_i | \mathbf{h}_i)$ from the training data?
Which smoothing technique will be used?
(Good-Turing, Jelinek-Mercer, Katz, Kneser-Ney,...)
Generalized Parallel Backoff etc.

Language Models – design decisions

1) How to factorize $P(w_1, w_2, \dots, w_m)$ into $\prod_{i=1..m} P(w_i | \mathbf{h}_i)$,
i.e. what conditions will be used as the context \mathbf{h}_i ?

**this
work**

2) What additional context information will be used
(apart from word forms),
e.g. stems, grammata, POS tags, word classes, ... ?

3) How to estimate $P(w_i | \mathbf{h}_i)$ from the training data?

Linear interpolation (what technique will be used?
Weights trained by EM (Baker, Mercer, Katz, Kneser-Ney, ...))

Generalized Parallel Backoff etc.

Language Models – design decisions

1) How to estimate $P(w_i | w_1, w_2, \dots, w_{i-1})$,
i.e. what context information will be used

**this
work**

$h_i = w_{i-n+1}, \dots, w_{i-1}$
(n-gram-based LMs)

2) What additional context information will be used
(apart from word forms),
e.g. stems, lemmata, POS tags, word classes, ... ?

3) How to estimate $P(w_i | h_i)$ from the training data?

Linear interpolation (Buckley, 1990)
Weights trained by EM (Moffitt, Mercer, Katz, 1996; Brill, 1996; ...)
Generalized Parallel Backoff etc.

**other
papers**

Outline

- Language Models (LM)
 - basics
 - design decisions
- **Post-ngram LM**
- Dependency LM
- Evaluation
- Conclusion & future plans

Post-ngram LM

In general: $P(s) = P(w_1, w_2, \dots, w_m) \approx \prod_{i=1..m} P(w_i | \mathbf{h}_i)$
 \mathbf{h}_i ... context (history) of word w_i

left-to-right factorization order

Bigram LM: $\mathbf{h}_i = w_{i-1}$ (one previous word)

Trigram LM: $\mathbf{h}_i = w_{i-2}, w_{i-1}$ (two previous words)

Post-ngram LM

In general: $P(s) = P(w_1, w_2, \dots, w_m) \approx \prod_{i=1..m} P(w_i | \mathbf{h}_i)$
 \mathbf{h}_i ... context (history) of word w_i

left-to-right factorization order

Bigram LM: $\mathbf{h}_i = w_{i-1}$ (one previous word)

Trigram LM: $\mathbf{h}_i = w_{i-2}, w_{i-1}$ (two previous words)

right-to-left factorization order

Post-bigram LM: $\mathbf{h}_i = w_{i+1}$ (one following word)

Post-trigram LM: $\mathbf{h}_i = w_{i+1}, w_{i+2}$ (two following words)

Post-ngram LM

In general: $P(s) = P(w_1, w_2, \dots, w_m) \approx \prod_{i=1..m} P(w_i | \mathbf{h}_i)$
 \mathbf{h}_i ... context (history) of word w_i

left-to-right factorization order

Bigram LM: $\mathbf{h}_i = w_{i-1}$ (one previous word)

Trigram LM: $\mathbf{h}_i = w_{i-2}, w_{i-1}$ (two previous words)

right-to-left factorization order

Post-bigram LM: $\mathbf{h}_i = w_{i+1}$ (one following word)

$P(\text{The dog barked again}) = P(\text{again} | \text{NONE}) \cdot P(\text{barked} | \text{again}) \cdot$
 $P(\text{dog} | \text{barked}) \cdot P(\text{The} | \text{dog})$

Outline

- Language Models (LM)
 - basics
 - design decisions
- Post-ngram LM
- **Dependency LM**
- Evaluation
- Conclusion & future plans

Dependency LM

- exploit the topology of dependency trees

The

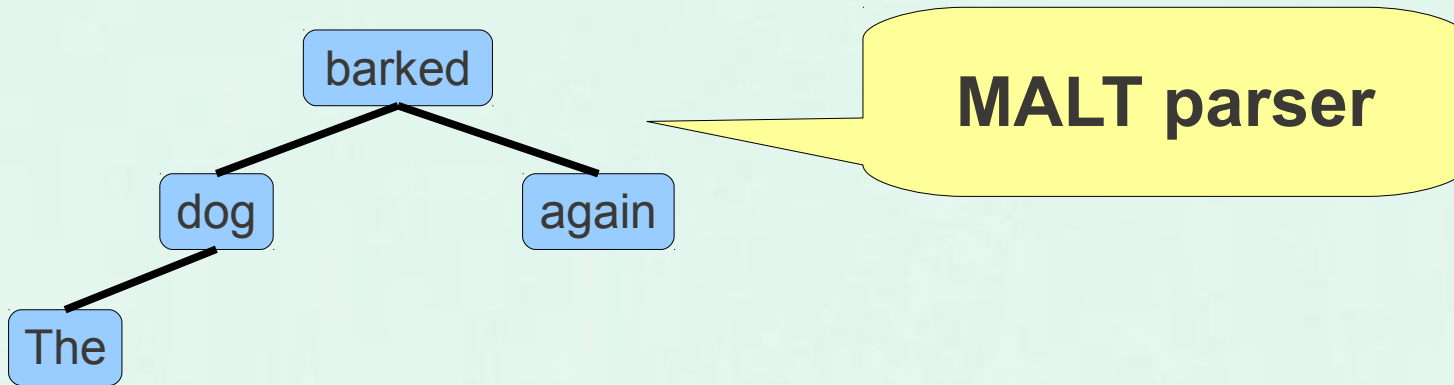
dog

barked

again

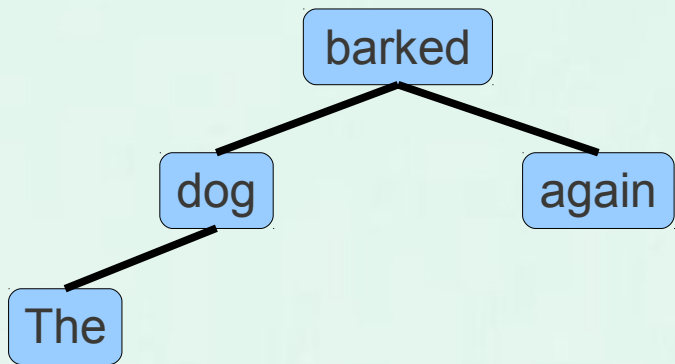
Dependency LM

- exploit the topology of dependency trees



Dependency LM

- exploit the topology of dependency trees



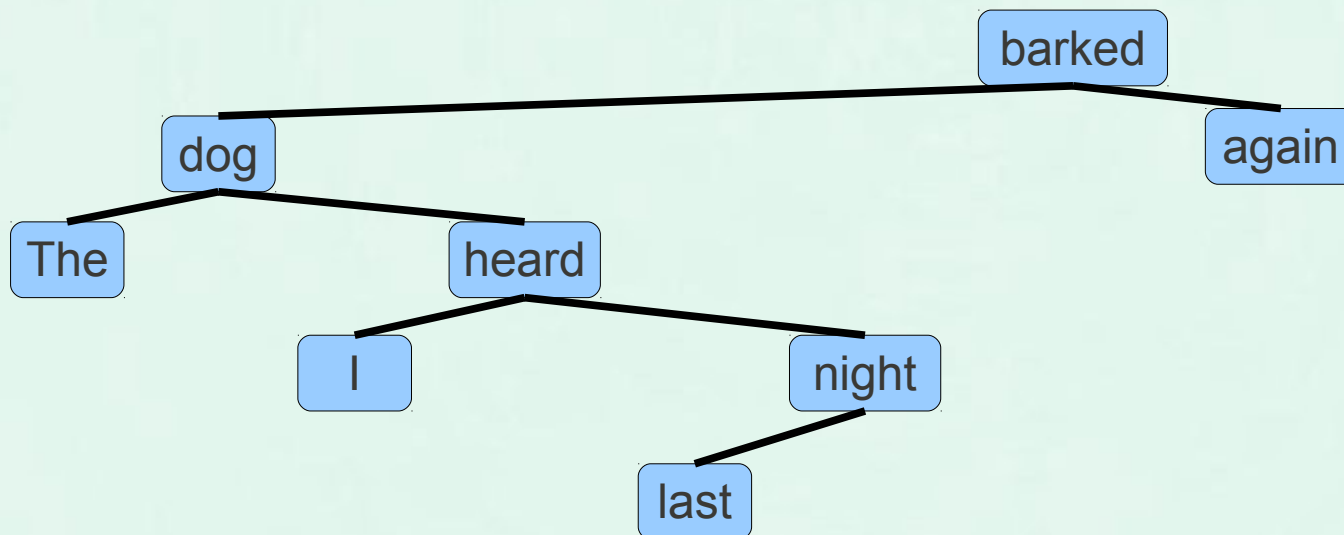
$$P(\text{The dog barked again}) = P(\text{The} \mid \text{dog}) \cdot P(\text{dog} \mid \text{barked}) \cdot P(\text{barked} \mid \text{NONE}) \cdot P(\text{again} \mid \text{barked})$$

$$h_i = \text{parent}(w_i)$$

Dependency LM

Long distance dependencies

The dog I heard last night barked again



Dependency LM

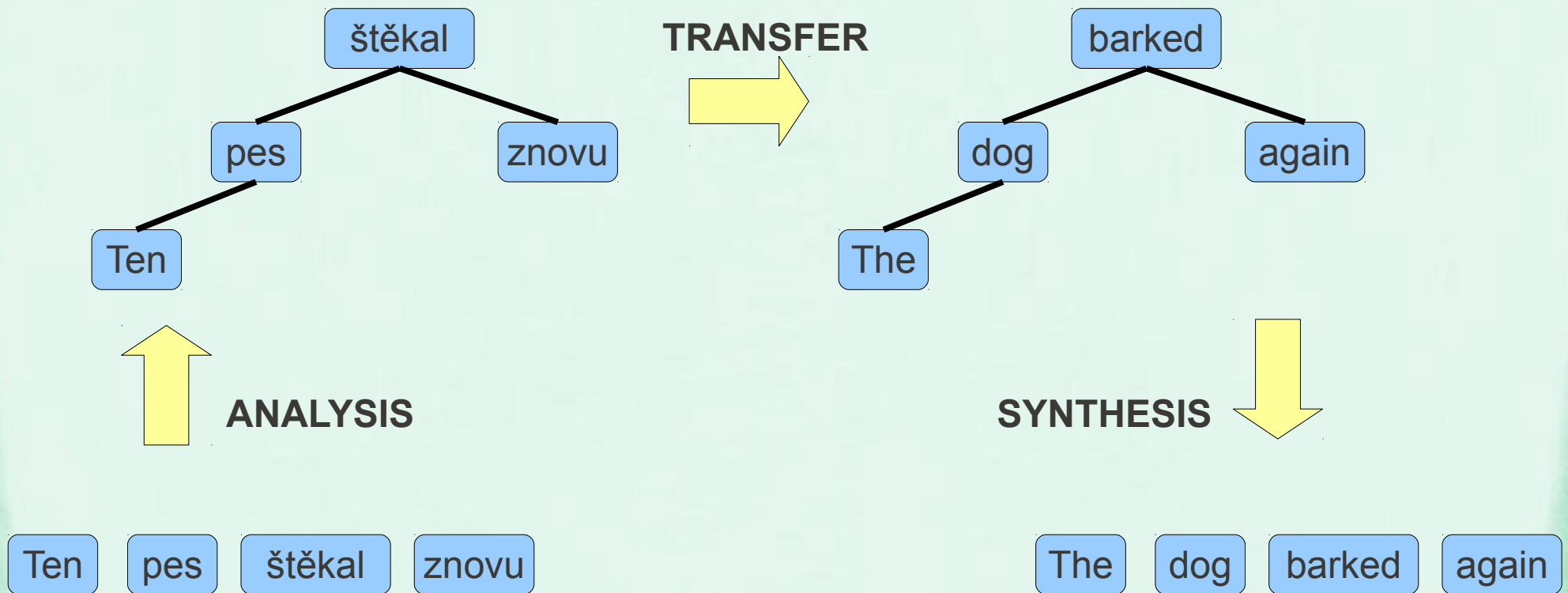
Motivation for usage

- How can we know the dependency structure without knowing the word-forms?

Dependency LM

Motivation for usage

- How can we know the dependency structure without knowing the word-forms?
- For example in **tree-to-tree machine translation**.

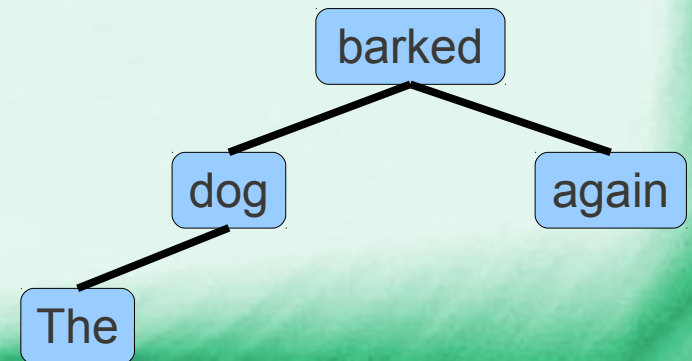


Dependency LM

Examples

- Model w_p
word form of parent

$$P(\text{The dog barked again}) =$$
$$P(\text{The} \mid \text{dog}) \cdot$$
$$P(\text{dog} \mid \text{barked}) \cdot$$
$$P(\text{barked} \mid \text{NONE}) \cdot$$
$$P(\text{again} \mid \text{barked})$$

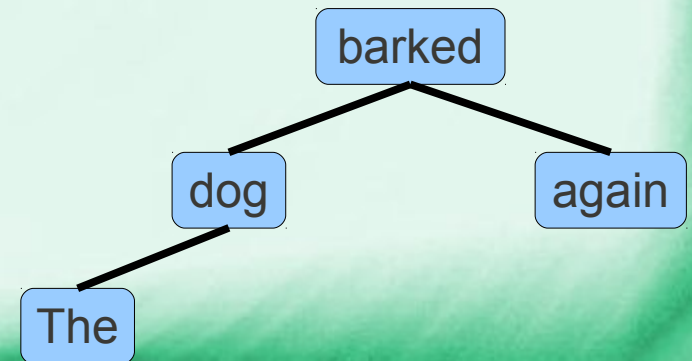


Dependency LM

Examples

- Model w_p, w_g
word form of parent, word form of grandparent

$$P(\text{The dog barked again}) =$$
$$P(\text{The} \mid \text{dog}, \text{barked}) \cdot$$
$$P(\text{dog} \mid \text{barked}, \text{NONE}) \cdot$$
$$P(\text{barked} \mid \text{NONE}, \text{NONE}) \cdot$$
$$P(\text{again} \mid \text{barked}, \text{NONE})$$

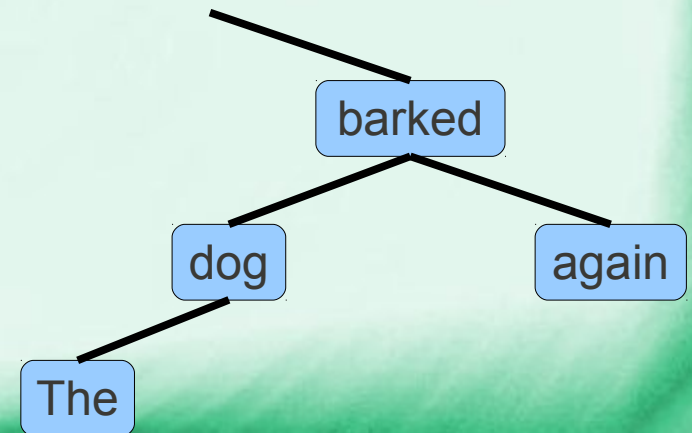


Dependency LM

Examples

- Model \mathbf{E}, \mathbf{wp}
edge direction, word form of parent

$$P(\text{The dog barked again}) =$$
$$P(\text{The} \mid \text{right}, \text{dog}) \cdot$$
$$P(\text{dog} \mid \text{right}, \text{barked}) \cdot$$
$$P(\text{barked} \mid \text{left}, \text{NONE}) \cdot$$
$$P(\text{again} \mid \text{left}, \text{barked})$$

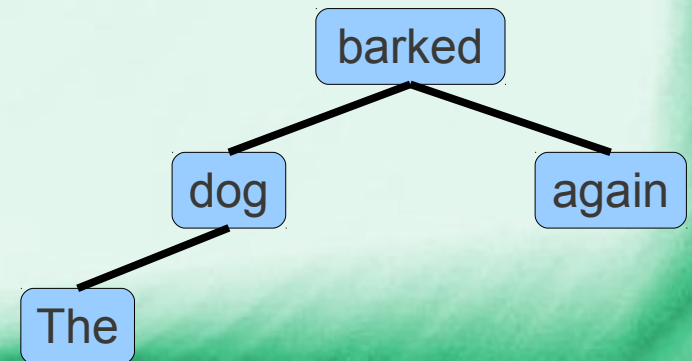


Dependency LM

Examples

- Model C, wp
number of children, word form of parent

$$P(\text{The dog barked again}) =$$
$$P(\text{The} \mid 0, \text{dog}) \cdot$$
$$P(\text{dog} \mid 1, \text{barked}) \cdot$$
$$P(\text{barked} \mid 2, \text{NONE}) \cdot$$
$$P(\text{again} \mid 0, \text{barked})$$



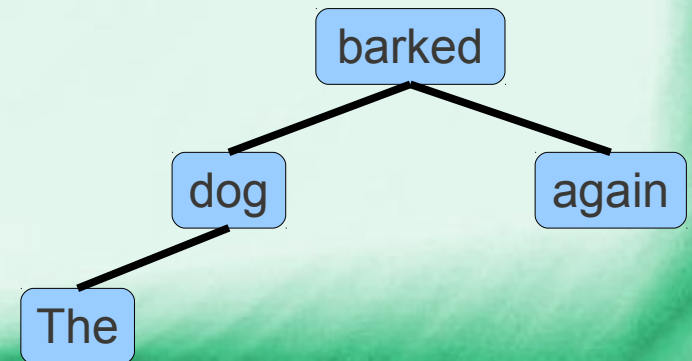
Dependency LM

Examples

- Model \mathbf{N}, \mathbf{wp}

the word is N^{th} child of its parent, word form of parent

$$\begin{aligned} P(\text{The dog barked again}) = & \\ P(\text{The} \mid 1, \text{dog}) \cdot & \\ P(\text{dog} \mid 1, \text{barked}) \cdot & \\ P(\text{barked} \mid 1, \text{NONE}) \cdot & \\ P(\text{again} \mid 2, \text{barked}) & \end{aligned}$$

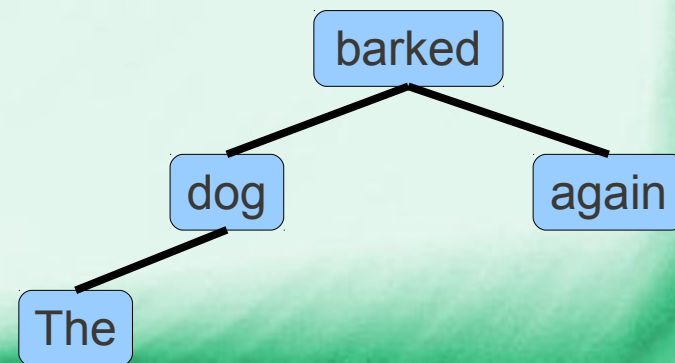


Dependency LM

Examples of additional context information

- Model t_p, w_p
POS tag of parent, word form of parent

$$P(\text{The dog barked again}) =$$
$$P(\text{The} \mid \text{NN}, \text{dog}) \cdot$$
$$P(\text{dog} \mid \text{VBD}, \text{barked}) \cdot$$
$$P(\text{barked} \mid \text{NONE}, \text{NONE}) \cdot$$
$$P(\text{again} \mid \text{VBD}, \text{barked})$$



Dependency LM

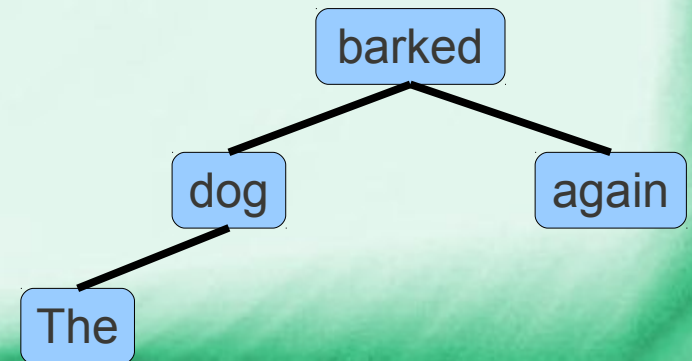
Examples of additional context information

- Model t_p, w_p

POS tag of parent, word form of parent

$$P(\text{The dog barked again}) =$$
$$P(\text{The} \mid \text{NN}, \text{dog}) \cdot$$
$$P(\text{dog} \mid \text{VBD}, \text{barked}) \cdot$$
$$P(\text{barked} \mid \text{NONE}, \text{NONE}) \cdot$$
$$P(\text{again} \mid \text{VBD}, \text{barked})$$

naïve tagger
assignes the most
frequent tag
for a given word

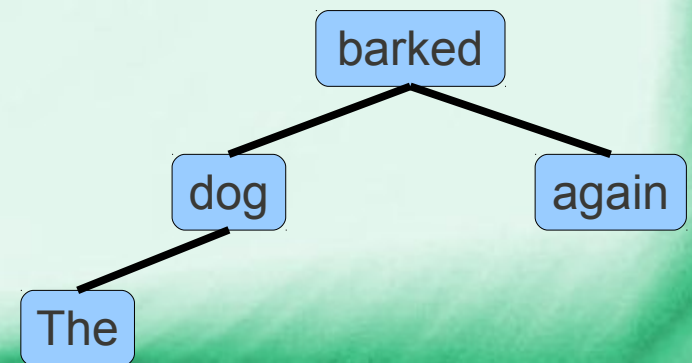


Dependency LM

Examples of additional context information

- Model T_p, w_p
coarse-grained POS tag of parent, word form of parent

$$P(\text{The dog barked again}) =$$
$$P(\text{The} \mid N, \text{dog}) \cdot$$
$$P(\text{dog} \mid V, \text{barked}) \cdot$$
$$P(\text{barked} \mid x, \text{NONE}) \cdot$$
$$P(\text{again} \mid V, \text{barked})$$

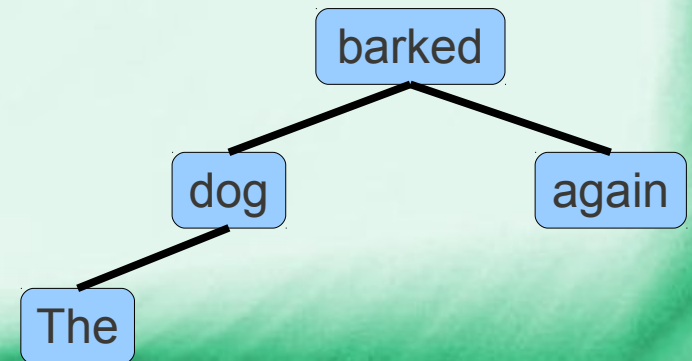


Dependency LM

Examples of additional context information

- Model $\mathbf{E}, \mathbf{C}, \mathbf{wp}, \mathbf{N}$
edge direction, # children, word form of parent,
word is N^{th} child of its parent

$$\begin{aligned} P(\text{The dog barked again}) = & \\ P(\text{The} \mid \text{right}, 0, \text{dog}, 1) \cdot & \\ P(\text{dog} \mid \text{right}, 1, \text{barked}, 1) \cdot & \\ P(\text{barked} \mid \text{left}, 2, \text{NONE}, 1) \cdot & \\ P(\text{again} \mid \text{left}, 0, \text{barked}, 2) & \end{aligned}$$



Outline

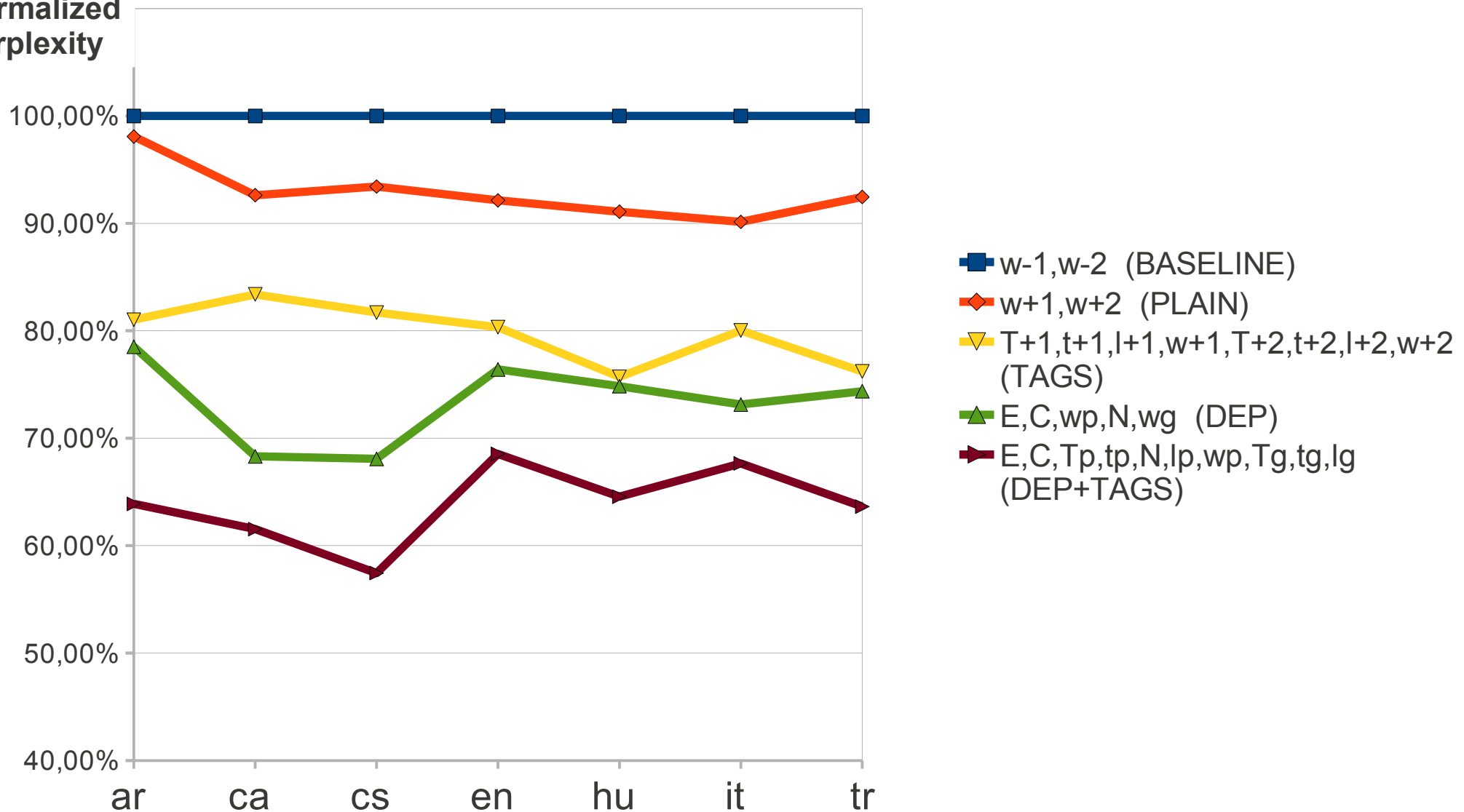
- Language Models (LM)
 - basics
 - design decisions
- Post-ngram LM
- Dependency LM
- **Evaluation**
- Conclusion & future plans

Evaluation

- Train and test data from CoNLL 2007 shared task
- 7 languages: Arabic, Catalan, Czech, English (450 000 tokens, 3 % OOV), Hungarian, Italian (75 000 tokens), and Turkish (26 % OOV)
- Cross-entropy = $-(1/|T|) \sum_{i=1..|T|} \log_2 P(w_i | \mathbf{h}_i)$,
measured on the test data T
- Perplexity = $2^{\text{Cross-entropy}}$
- Lower perplexity ~ better LM
- Baseline ... trigram LM
- 4 experimental settings: PLAIN, TAGS, DEP, DEP+TAGS

Evaluation

normalized
perplexity



Conclusion

Findings confirmed
for all seven languages

Improvement over baseline
for English

- Post-trigram better than trigram **PLAIN 8 %**
- Post-bigram better than bigram
- Additional context (POS & lemma) helps **TAGS 20 %**
- Dependency structure helps even more **DEP 24 %**
- The best perplexity achieved with **DEP+TAGS 31 %**

Future plans

- Investigate the reason for better post-ngram LM perplexity
- Extrinsic evaluation
 - Post-ngram LM in speech recognition
 - Dependency LM in tree-to-tree machine translation
- Better smoothing using Generalized Parallel Backoff
- Bigger LM for real applications

Thank you