

Statistical Machine Translation with Weighted Grammars

Matthias BÜchse

Chair of Foundations of Programming
Institute of Theoretical Computer Science
Technische Universität Dresden

June 15, 2011

Outline

Statistical Machine Translation

Weighted Grammars as a Feature

Statistical Machine Translation with Weighted Grammars

Ambiguity in Natural Language

Example (Chiang 2007)

- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

Ambiguity in Natural Language

Example (Chiang 2007)

- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

Is e a correct translation of f?

Ambiguity in Natural Language

Example (Chiang 2007)

- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

Is e a correct translation of f ?



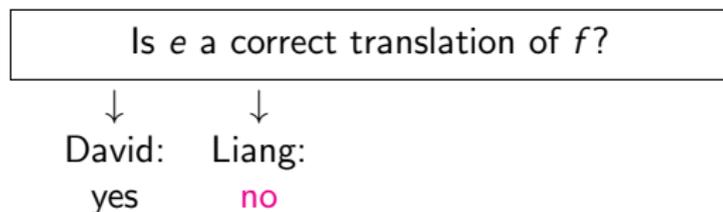
David:

yes

Ambiguity in Natural Language

Example (Chiang 2007)

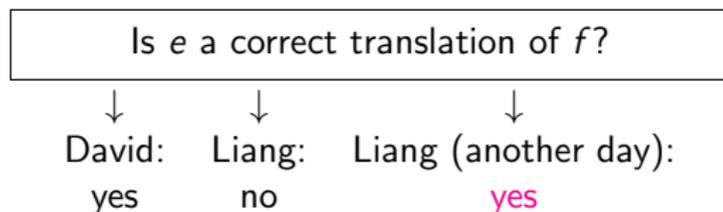
- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”



Ambiguity in Natural Language

Example (Chiang 2007)

- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”



Ambiguity in Natural Language

Example (Chiang 2007)

- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

What is $P(e | f)$?

Ambiguity in Natural Language

Example (Chiang 2007)

- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

What is $P(e | f)$?

↓
.42

Ambiguity in Natural Language

Example (Chiang 2007)

- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

What is $P(e | f)$?

↓
.42

frequentist limit of relative frequencies

Bayesian degree of belief (subjective)

Search (aka Decoding) in SMT

From: German ▼  To: English ▼ Translate

Ich möchte diesen Teppich nicht kaufen.



 Listen

German to English translation

I would not buy this carpet.

 Listen

New! Click the words above to view alternate translations. [Dismiss](#)

given $p \in \mathcal{M}(E | F)$

Search (aka Decoding) in SMT

From: German ▼  To: English ▼ Translate

Ich möchte diesen Teppich nicht kaufen.



 Listen

German to English translation

I would not buy this carpet.

 Listen

New! Click the words above to view alternate translations. [Dismiss](#)

given $p \in \mathcal{M}(E | F)$ ($\sum_e p(e | f) = 1$ for every f)

Search (aka Decoding) in SMT

From: German ▼  To: English ▼ Translate

Ich möchte diesen Teppich nicht kaufen.



 Listen

German to English translation

I would not buy this carpet.

 Listen

New! Click the words above to view alternate translations. [Dismiss](#)

given $p \in \mathcal{M}(E | F)$ and $f \in F$,

Search (aka Decoding) in SMT

From: German ▼  To: English ▼ Translate

Ich möchte diesen Teppich nicht kaufen.



 Listen

German to English translation

I would not buy this carpet.

 Listen

New! Click the words above to view alternate translations. [Dismiss](#)

given $p \in \mathcal{M}(E | F)$ and $f \in F$,
determine \hat{e} such that for every $e \in E$:

$$p(\hat{e} | f) \geq p(e | f) . \quad (\text{not unique})$$

Search (aka Decoding) in SMT

From: German ▼ To: English ▼ Translate

Ich möchte diesen Teppich nicht kaufen.



Listen

German to English translation

I would not buy this carpet.

Listen

New! Click the words above to view alternate translations. [Dismiss](#)

given $p \in \mathcal{M}(E | F)$ and $f \in F$,
determine \hat{e} such that for every $e \in E$:

$$p(\hat{e} | f) \geq p(e | f) . \quad (\text{not unique})$$

for short

$$\hat{e} = \operatorname{argmax}_{e \in E} p(e | f) .$$

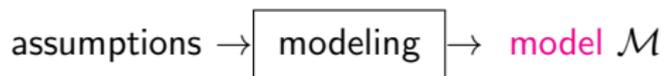
Building an SMT System

assumptions

training data

test data

Building an SMT System

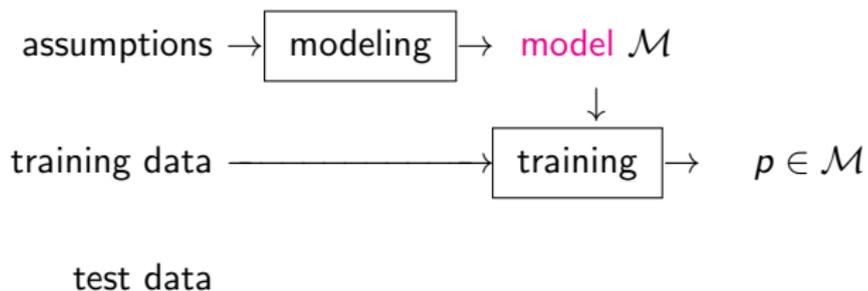


training data

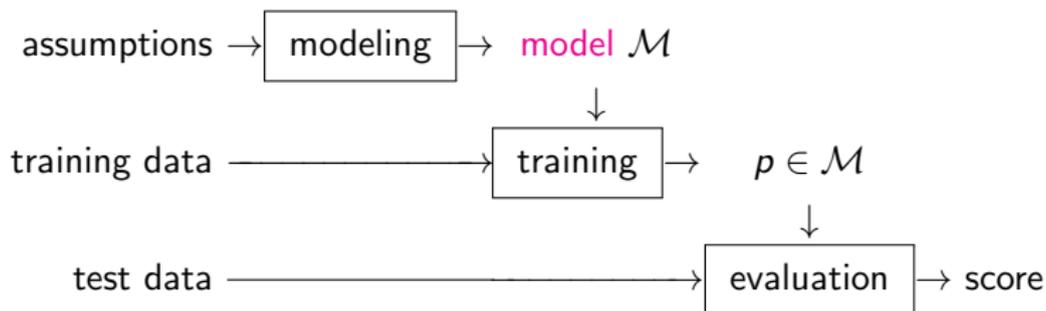
test data

“Every time I fire a linguist, recognition rate goes up”
(attributed to Frederick Jelinek)

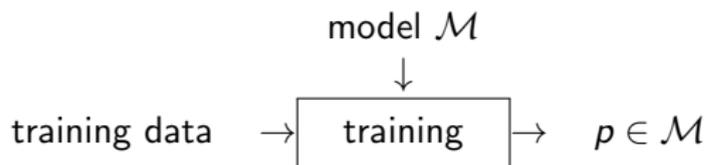
Building an SMT System



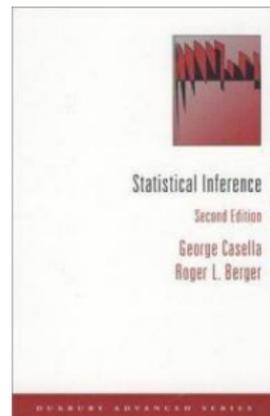
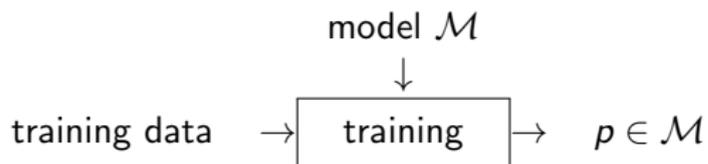
Building an SMT System



Training = Point Estimation

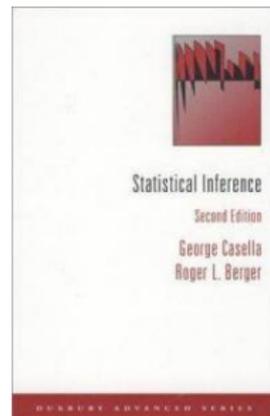
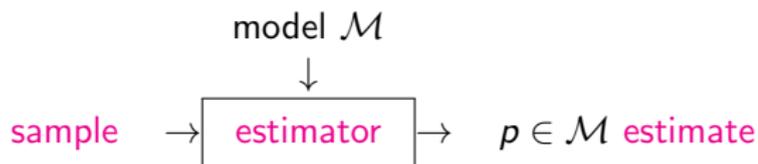


Training = Point Estimation



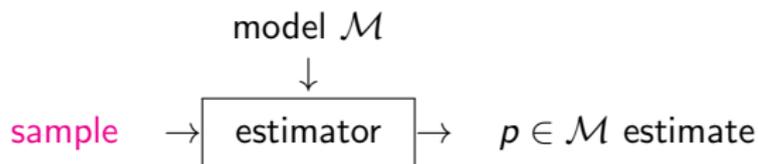
Statistical Inference, Ch. 7:
Point Estimation
(Casella and Berger 2002)

Training = Point Estimation

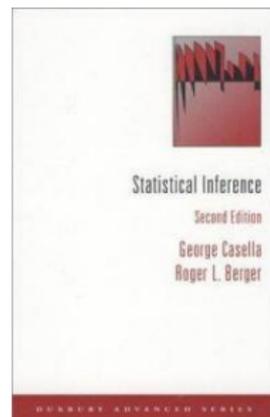


Statistical Inference, Ch. 7:
Point Estimation
(Casella and Berger 2002)

Training = Point Estimation

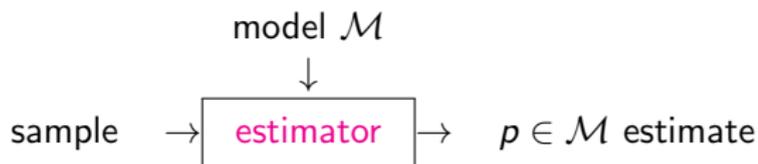


- sample
- ▶ representative of all translation situations
 - ▶ e. g., Hong Kong Hansards (parliament proceedings)



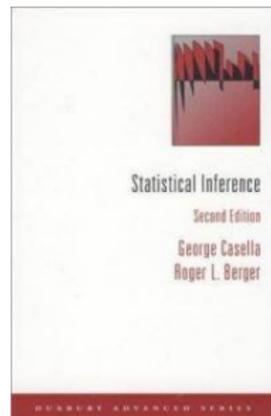
Statistical Inference, Ch. 7:
Point Estimation
(Casella and Berger 2002)

Training = Point Estimation



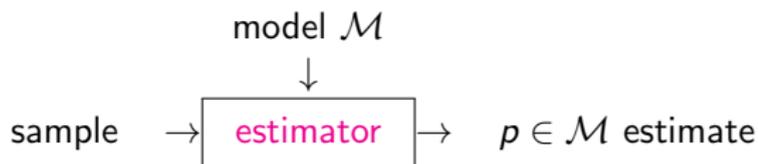
- sample
- ▶ representative of all translation situations
 - ▶ e. g., Hong Kong Hansards (parliament proceedings)

- estimators
- ▶ maximum likelihood
 - ▶ maximum a-posteriori
 - ▶ minimum risk
 - ▶ maximum entropy



Statistical Inference, Ch. 7:
Point Estimation
(Casella and Berger 2002)

Training = Point Estimation



sample ▶ representative of all translation situations

▶ e. g., Hong Kong Hansards (parliament proceedings)

estimators ▶ maximum likelihood

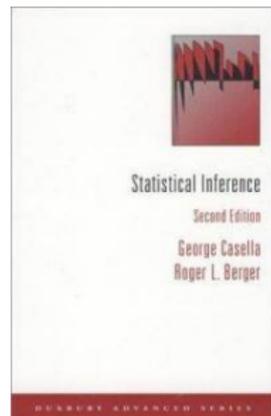
▶ maximum a-posteriori

▶ minimum risk

▶ maximum entropy

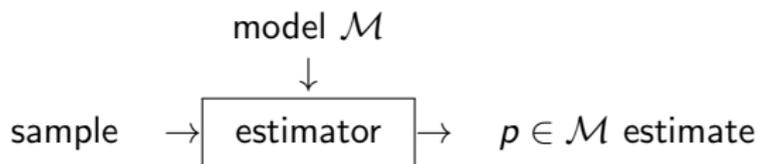
properties ▶ loss function optimality (Statistical Decision Theory)

▶ asymptotic optimality (consistency)



Statistical Inference, Ch. 7:
Point Estimation
(Casella and Berger 2002)

Training = Point Estimation



sample ▶ representative of all translation situations

▶ e. g., Hong Kong Hansards (parliament proceedings)

estimators ▶ maximum likelihood

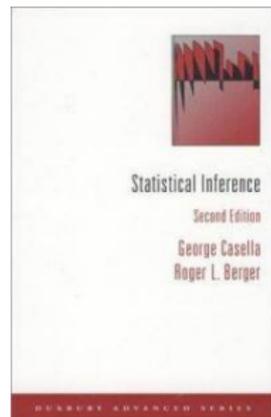
▶ maximum a-posteriori

▶ minimum risk

▶ maximum entropy

properties ▶ loss function optimality (Statistical Decision Theory)

▶ asymptotic optimality (consistency)



Statistical Inference, Ch. 7:
Point Estimation
(Casella and Berger 2002)

Sample?

- ▶ data sparsity

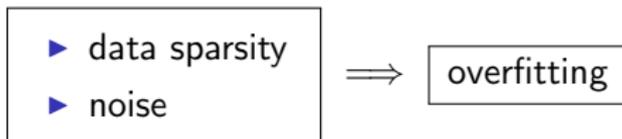
Sample?

- ▶ data sparsity
- ▶ noise

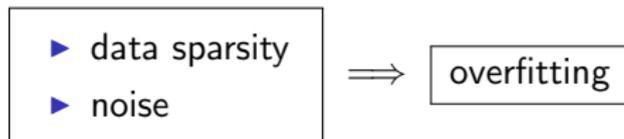
Sample?

- ▶ data sparsity
- ▶ noise

Sample?



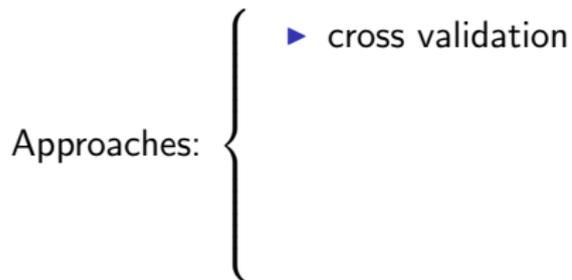
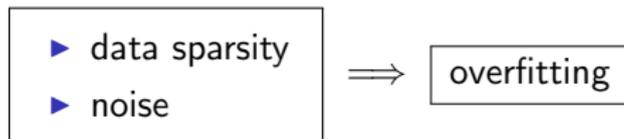
Sample?



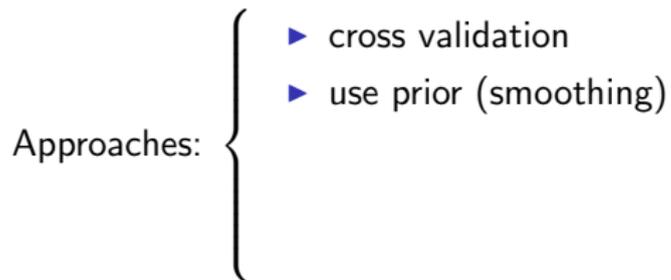
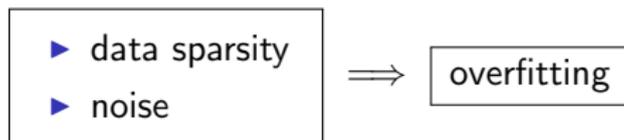
Approaches:



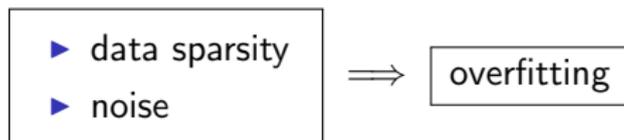
Sample?



Sample?

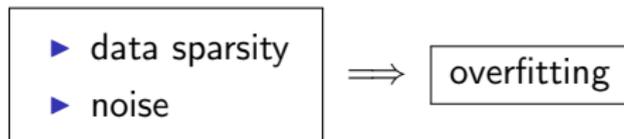


Sample?



- Approaches: {
- ▶ cross validation
 - ▶ use prior (smoothing)
 - ▶ use ad-hoc estimator (heuristics)

Sample?



- Approaches: {
- ▶ cross validation
 - ▶ use prior (smoothing)
 - ▶ use ad-hoc estimator (heuristics)
 - ▶ adjust the model

Adjusting the Model

Source-Channel Model (Brown et al. 1993)

- ▶ $\mathcal{M}_{\text{LM}} \subseteq \mathcal{M}(E)$ language model
- ▶ $\mathcal{M}_{\text{TM}} \subseteq \mathcal{M}(F | E)$ translation model

Adjusting the Model

Source-Channel Model (Brown et al. 1993)

- ▶ $\mathcal{M}_{\text{LM}} \subseteq \mathcal{M}(E)$ language model
- ▶ $\mathcal{M}_{\text{TM}} \subseteq \mathcal{M}(F | E)$ translation model
- ▶ then

$$\mathcal{M} = \{p \in \mathcal{M}(E | F) \mid \exists p_1 \in \mathcal{M}_{\text{LM}}, p_2 \in \mathcal{M}_{\text{TM}}:$$

$$p(e | f) = \frac{p_1(e) \cdot p_2(f | e)}{\sum_{e' \in E} p_1(e') \cdot p_2(f | e')}\}$$

Adjusting the Model

Source-Channel Model (Brown et al. 1993)

- ▶ $\mathcal{M}_{\text{LM}} \subseteq \mathcal{M}(E)$ language model
- ▶ $\mathcal{M}_{\text{TM}} \subseteq \mathcal{M}(F | E)$ translation model
- ▶ then

$$\mathcal{M} = \left\{ p \in \mathcal{M}(E | F) \mid \exists p_1 \in \mathcal{M}_{\text{LM}}, p_2 \in \mathcal{M}_{\text{TM}}: \right. \\ \left. p(e | f) = \frac{p_1(e) \cdot p_2(f | e)}{\sum_{e' \in E} p_1(e') \cdot p_2(f | e')} \right\}$$

- ▶ estimate p_1 and p_2 independently on appropriate data

Adjusting the Model

Log-linear Model (Berger, Della Pietra, and Della Pietra 1996; Och and Ney 2002)

- ▶ $h_1, \dots, h_n: E \times F \rightarrow \mathbb{R}$ features

Adjusting the Model

Log-linear Model (Berger, Della Pietra, and Della Pietra 1996; Och and Ney 2002)

- ▶ $h_1, \dots, h_n: E \times F \rightarrow \mathbb{R}$ features
- ▶ consider

$$\mathcal{M} = \{p \in \mathcal{M}(E | F) \mid \forall i: \mathbb{E}_p[h_i] = \mathbb{E}_{\tilde{p}}[h_i]\}$$

- ▶ (and maximize entropy)

Adjusting the Model

Log-linear Model (Berger, Della Pietra, and Della Pietra 1996; Och and Ney 2002)

- ▶ $h_1, \dots, h_n: E \times F \rightarrow \mathbb{R}$ features
- ▶ consider

$$\mathcal{M} = \{p \in \mathcal{M}(E | F) \mid \forall i: \mathbb{E}_p[h_i] = \mathbb{E}_{\bar{p}}[h_i]\}$$

- ▶ (and maximize entropy)
- ▶ or consider

$$\mathcal{M} = \{p \in \mathcal{M}(E | F) \mid \exists \lambda_1, \dots, \lambda_n \in \mathbb{R}:$$

$$p(e | f) = \frac{\exp \sum_i \lambda_i \cdot h_i(e, f)}{\sum_{e' \in E} \exp \sum_i \lambda_i \cdot h_i(e', f)}\}$$

- ▶ (and maximize likelihood)

Adjusting the Model

Log-linear Model Example

- ▶ $h_1(e, f) = \log p_1(e)$,
- ▶ $h_2(e, f) = \log p_2(f | e)$,
- ▶ $\lambda_1 = 1$,
- ▶ $\lambda_2 = 1$,
- ▶ then

$$\begin{aligned} p(e | f) &= \frac{\exp(\lambda_1 \cdot h_1(e, f) + \lambda_2 \cdot h_2(e, f))}{\sum_{e' \in E} \exp(\lambda_1 \cdot h_1(e', f) + \lambda_2 \cdot h_2(e', f))} \\ &= \frac{\exp(\log p_1(e) + \log p_2(f | e))}{\sum_{e' \in E} \exp(\log p_1(e') + \log p_2(f | e))} \\ &= \frac{p_1(e) \cdot p_2(f | e)}{\sum_{e' \in E} p_1(e') \cdot p_2(f | e)} \end{aligned}$$

Estimating Conditional Models

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ define

$$\mathcal{M}' = \{p' \in \mathcal{M}(E, F) \mid \exists p \in \mathcal{M}, p_F \in \mathcal{M}(F): \\ p'(e, f) = p(e | f) \cdot p_F(f)\}$$

Estimating Conditional Models

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ define

$$\mathcal{M}' = \{p' \in \mathcal{M}(E, F) \mid \exists p \in \mathcal{M}, p_F \in \mathcal{M}(F): \\ p'(e, f) = p(e | f) \cdot p_F(f)\}$$

- ▶ bijection: $p' \leftrightarrow (p, p_F)$ iff $p'(e, f) = p(e | f) \cdot p_F(f)$
- ▶ if $P(e, f) = p'(e, f)$, then $P(e | f) = p(e | f)$,
- ▶ but p and p_F are completely independent,

Estimating Conditional Models

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ define

$$\mathcal{M}' = \{p' \in \mathcal{M}(E, F) \mid \exists p \in \mathcal{M}, p_F \in \mathcal{M}(F): \\ p'(e, f) = p(e | f) \cdot p_F(f)\}$$

- ▶ bijection: $p' \leftrightarrow (p, p_F)$ iff $p'(e, f) = p(e | f) \cdot p_F(f)$
 - ▶ if $P(e, f) = p'(e, f)$, then $P(e | f) = p(e | f)$,
 - ▶ but p and p_F are completely independent,
- ↪ do estimation on \mathcal{M}' , throw p_F away

More Recent Estimators

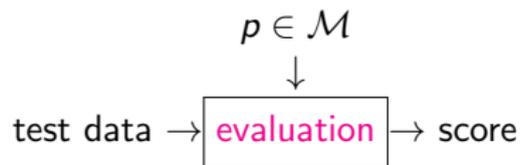
- ▶ minimum-error-rate training (Och 2003)
- ▶ minimum-risk annealing (Smith and Eisner 2006)
- ▶ large-margin methods

Evaluation

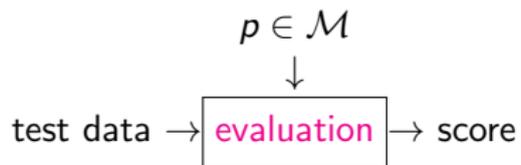


(due to Saša Hasan)

Evaluation



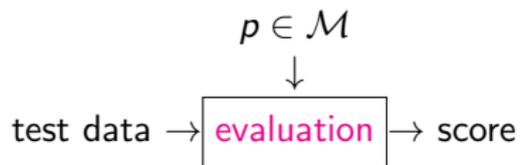
Evaluation



training data \cap test data = \emptyset



Evaluation

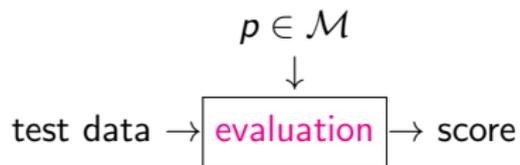


$$\text{training data} \cap \text{test data} = \emptyset$$

method	how it works	goal
BLEU (Papineni et al. 2002)	n -gram precision, brevity penalty	higher values



Evaluation

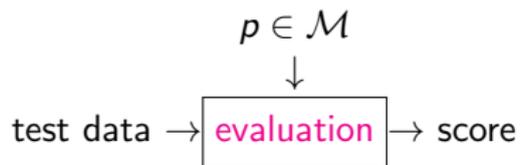


$$\text{training data} \cap \text{test data} = \emptyset$$

method	how it works	goal
BLEU (Papineni et al. 2002)	n -gram precision, brevity penalty	higher values
TER (Snover et al. 2006)	edit distance	lower values



Evaluation



$$\text{training data} \cap \text{test data} = \emptyset$$

method	how it works	goal
BLEU (Papineni et al. 2002)	n -gram precision, brevity penalty	higher values
TER (Snover et al. 2006)	edit distance	lower values
Meteor (Banerjee and Lavie 2005)	alignments, unigram precision, unigram recall, penalty	higher values



Outline

Statistical Machine Translation

Weighted Grammars as a Feature

Statistical Machine Translation with Weighted Grammars

Example (Chiang 2007)

Recall

- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

Example (Chiang 2007)

Recall

- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

Synchronous Context-Free Grammar G :

$\pi_1: S \rightarrow \langle S X, S X \rangle$

$\pi_2: S \rightarrow \langle X, X \rangle$

$\pi_3: X \rightarrow \langle \text{yu } X_{[1]} \text{ you } X_{[2]}, \text{ have } X_{[2]} \text{ with } X_{[1]} \rangle$

$\pi_4: X \rightarrow \langle X_{[1]} \text{ de } X_{[2]}, \text{ the } X_{[2]} \text{ that } X_{[1]} \rangle$

$\pi_5: X \rightarrow \langle X \text{ zhiyi, one of } X \rangle$

$\pi_6: X \rightarrow \langle \text{Aozhou, Australia} \rangle$

$\pi_7: X \rightarrow \langle \text{Beihan, North Korea} \rangle$

$\pi_8: X \rightarrow \langle \text{shi, is} \rangle$

$\pi_9: X \rightarrow \langle \text{bangjiao, diplomatic relations} \rangle$

$\pi_{10}: X \rightarrow \langle \text{shaoshu guojia, few countries} \rangle$

Derivation

$\langle S, S \rangle$

Derivation

$$\begin{array}{l} \langle S, S \rangle \\ \xRightarrow{\pi_1} \langle S X, S X \rangle \end{array}$$

Derivation

$$\begin{array}{l} \langle S, S \rangle \\ \xRightarrow{\pi_1} \langle S X, S X \rangle \end{array}$$

Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \end{aligned}$$

Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{\boxed{1}} X_{\boxed{2}}, S X_{\boxed{1}} X_{\boxed{2}} \rangle \end{aligned}$$

Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \end{aligned}$$

Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \end{aligned}$$

Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_6} & \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle \end{aligned}$$

Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_6} & \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle \end{aligned}$$

Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xRightarrow{\pi_1} & \langle S X, S X \rangle \\ \xRightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_6} & \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_8} & \langle \text{Aozhou } \text{shi } X, \text{Australia } \text{is } X \rangle \end{aligned}$$

Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xRightarrow{\pi_1} & \langle S X, S X \rangle \\ \xRightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_6} & \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_8} & \langle \text{Aozhou shi } X, \text{Australia is } X \rangle \end{aligned}$$

Derivation

$\langle S, S \rangle$
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi}, \text{Australia is one of } X \rangle$

Derivation

$\langle S, S \rangle$
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$

Derivation

$\langle S, S \rangle$
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$
 $\text{Australia is one of the } X_{[2]} \text{ that } X_{[1]} \rangle$

Derivation

$\langle S, S \rangle$
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$
 $\quad \text{Australia is one of the } X_{[2]} \text{ that } X_{[1]} \rangle$

Derivation

$\langle S, S \rangle$
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$
 Australia is one of the $X_{[2]}$ that $X_{[1]}$ \rangle
 $\xRightarrow{\pi_3} \langle \text{Aozhou shi } \text{yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
 Australia is one of the $X_{[2]}$ that have $X_{[0]}$ with $X_{[1]}$ \rangle

Derivation

$\langle S, S \rangle$
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$
 Australia is one of the $X_{[2]}$ that $X_{[1]}$ \rangle
 $\xRightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
 Australia is one of the $X_{[2]}$ that have $X_{[0]}$ with $X_{[1]}$ \rangle

Derivation

$\langle S, S \rangle$
 $\xrightarrow{\pi_1} \langle S X, S X \rangle$
 $\xrightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
 $\xrightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
 $\xrightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
 $\xrightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
 $\xrightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$
 $\xrightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$
 Australia is one of the } X_{[2]} \text{ that } X_{[1]} \rangle
 $\xrightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
 Australia is one of the } X_{[2]} \text{ that have } X_{[0]} \text{ with } X_{[1]} \rangle
 $\xrightarrow{\pi_7} \langle \text{Aozhou shi yu } \text{Beihan} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
 Australia is one of the } X_{[2]} \text{ that have } X_{[0]} \text{ with } \text{North Korea} \rangle

Derivation

$\langle S, S \rangle$
 $\xrightarrow{\pi_1} \langle S X, S X \rangle$
 $\xrightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
 $\xrightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
 $\xrightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
 $\xrightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
 $\xrightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$
 $\xrightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$
 Australia is one of the } X_{[2]} \text{ that } X_{[1]} \rangle
 $\xrightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
 Australia is one of the } X_{[2]} \text{ that have } X_{[0]} \text{ with } X_{[1]} \rangle
 $\xrightarrow{\pi_7} \langle \text{Aozhou shi yu Beihan you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
 Australia is one of the } X_{[2]} \text{ that have } X_{[0]} \text{ with North Korea} \rangle

Derivation

- $\langle S, S \rangle$
 $\xrightarrow{\pi_1} \langle S X, S X \rangle$
 $\xrightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
 $\xrightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
 $\xrightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
 $\xrightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
 $\xrightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$
 $\xrightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$
Australia is one of the $X_{[2]}$ that $X_{[1]}$ \rangle
 $\xrightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
Australia is one of the $X_{[2]}$ that have $X_{[0]}$ with $X_{[1]}$ \rangle
 $\xrightarrow{\pi_7} \langle \text{Aozhou shi yu Beihan you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
Australia is one of the $X_{[2]}$ that have $X_{[0]}$ with North Korea \rangle
 $\xrightarrow{\pi_9} \langle \text{Aozhou shi yu Beihan you } \text{bangjiao} \text{ de } X_{[2]} \text{ zhiyi,}$
Australia is one of the $X_{[2]}$ that have **diplomatic relations** with North Korea \rangle

Derivation

$\langle S, S \rangle$

$\xRightarrow{\pi_1} \langle S X, S X \rangle$

$\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$

$\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$

$\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$

$\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$

$\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$

$\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$
 $\text{Australia is one of the } X_{[2]} \text{ that } X_{[1]} \rangle$

$\xRightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
 $\text{Australia is one of the } X_{[2]} \text{ that have } X_{[0]} \text{ with } X_{[1]} \rangle$

$\xRightarrow{\pi_7} \langle \text{Aozhou shi yu Beihan you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
 $\text{Australia is one of the } X_{[2]} \text{ that have } X_{[0]} \text{ with North Korea} \rangle$

$\xRightarrow{\pi_9} \langle \text{Aozhou shi yu Beihan you bangjiao de } X_{[2]} \text{ zhiyi,}$
 $\text{Australia is one of the } X_{[2]} \text{ that have diplomatic relations with North Korea} \rangle$

Derivation

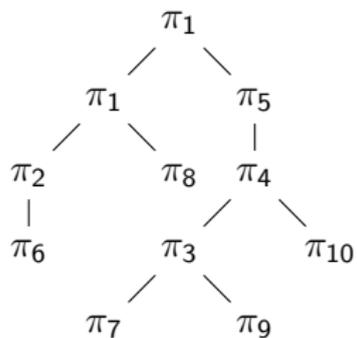
- $\langle S, S \rangle$
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$
Australia is one of the $X_{[2]}$ that $X_{[1]}$ \rangle
 $\xRightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
Australia is one of the $X_{[2]}$ that have $X_{[0]}$ with $X_{[1]}$ \rangle
 $\xRightarrow{\pi_7} \langle \text{Aozhou shi yu Beihan you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
Australia is one of the $X_{[2]}$ that have $X_{[0]}$ with North Korea \rangle
 $\xRightarrow{\pi_9} \langle \text{Aozhou shi yu Beihan you bangjiao de } X_{[2]} \text{ zhiyi,}$
Australia is one of the $X_{[2]}$ that have diplomatic relations with North Korea \rangle
 $\xRightarrow{\pi_{10}} \langle \text{Aozhou shi yu Beihan you bangjiao de } \text{shaoshu guojia} \text{ zhiyi,}$
Australia is one of the few countries that have diplomatic relations with N. K. \rangle

Derivation

- $\langle S, S \rangle$
- $\xRightarrow{\pi_1} \langle S X, S X \rangle$
- $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
- $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
- $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
- $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
- $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$
- $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$
 $\text{Australia is one of the } X_{[2]} \text{ that } X_{[1]} \rangle$
- $\xRightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
 $\text{Australia is one of the } X_{[2]} \text{ that have } X_{[0]} \text{ with } X_{[1]} \rangle$
- $\xRightarrow{\pi_7} \langle \text{Aozhou shi yu Beihan you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$
 $\text{Australia is one of the } X_{[2]} \text{ that have } X_{[0]} \text{ with North Korea} \rangle$
- $\xRightarrow{\pi_9} \langle \text{Aozhou shi yu Beihan you bangjiao de } X_{[2]} \text{ zhiyi,}$
 $\text{Australia is one of the } X_{[2]} \text{ that have diplomatic relations with North Korea} \rangle$
- $\xRightarrow{\pi_{10}} \langle \text{Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi,}$
 $\text{Australia is one of the few countries that have diplomatic relations with N. K.} \rangle$

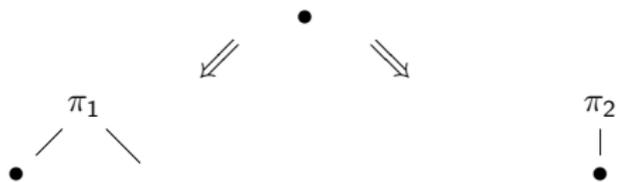
Inducing a Model

Derivation Tree



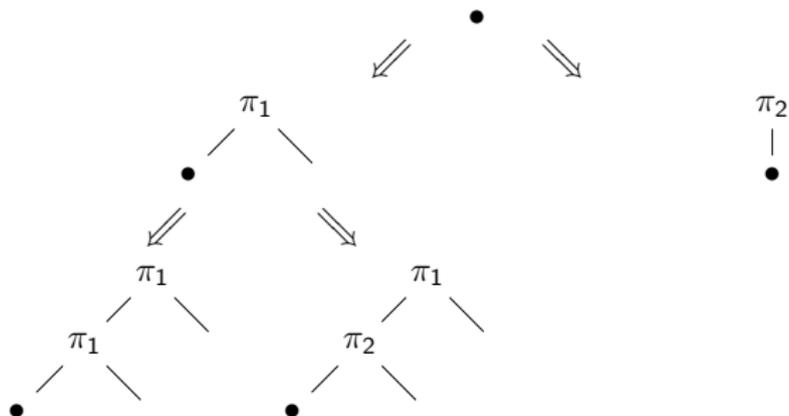
Inducing a Model

Decision Tree



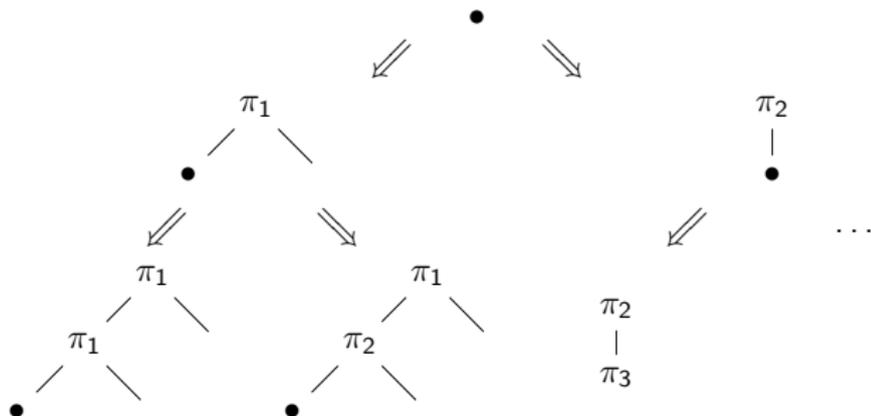
Inducing a Model

Decision Tree



Inducing a Model

Decision Tree



Inducing a Model

Assignments

grammar G } “distribution” $p'(e, d, f)$
assignment p' }

model $\mathcal{M}_G = \{p \in \mathcal{M}(E, F) \mid \exists p' : p(e, f) = \sum_d p'(e, d, f)\}$

Proper and Consistent Assignments

Let $q \in [0, 1]$.

$$\begin{aligned}\pi_1: S &\rightarrow \langle S S, S S \rangle \# q \\ \pi_2: S &\rightarrow \langle \varepsilon, \varepsilon \rangle \quad \# 1 - q\end{aligned}$$

Proper and Consistent Assignments

Let $q \in [0, 1]$.

$$\pi_1: S \rightarrow \langle S S, S S \rangle \# q$$

$$\pi_2: S \rightarrow \langle \varepsilon, \varepsilon \rangle \# 1 - q$$

Consider $\sum_d p(\varepsilon, d, \varepsilon)$. This is the least solution to

$$Z = q \cdot Z \cdot Z + (1 - q) .$$

Proper and Consistent Assignments

Let $q \in [0, 1]$.

$$\pi_1: S \rightarrow \langle S S, S S \rangle \# q$$

$$\pi_2: S \rightarrow \langle \varepsilon, \varepsilon \rangle \# 1 - q$$

Consider $\sum_d p(\varepsilon, d, \varepsilon)$. This is the least solution to

$$Z = q \cdot Z \cdot Z + (1 - q) \cdot$$

Reordering gives ($q \neq 0$)

$$0 = Z^2 - \frac{1}{q} \cdot Z + \frac{1 - q}{q}$$

Proper and Consistent Assignments

Let $q \in [0, 1]$.

$$\begin{aligned}\pi_1: S &\rightarrow \langle S S, S S \rangle \# q \\ \pi_2: S &\rightarrow \langle \varepsilon, \varepsilon \rangle \quad \# 1 - q\end{aligned}$$

Consider $\sum_d p(\varepsilon, d, \varepsilon)$. This is the least solution to

$$Z = q \cdot Z \cdot Z + (1 - q).$$

Reordering gives ($q \neq 0$)

$$0 = Z^2 - \frac{1}{q} \cdot Z + \frac{1 - q}{q}$$

Then

$$Z_1 = 1 \qquad Z_2 = \frac{1}{q} - 1$$

Proper and Consistent Assignments

Let $q \in [0, 1]$.

$$\begin{aligned}\pi_1: S &\rightarrow \langle S S, S S \rangle \# q \\ \pi_2: S &\rightarrow \langle \varepsilon, \varepsilon \rangle \quad \# 1 - q\end{aligned}$$

Consider $\sum_d p(\varepsilon, d, \varepsilon)$. This is the least solution to

$$Z = q \cdot Z \cdot Z + (1 - q).$$

Reordering gives ($q \neq 0$)

$$0 = Z^2 - \frac{1}{q} \cdot Z + \frac{1 - q}{q}$$

Then

$$Z_1 = 1 \qquad Z_2 = \frac{1}{q} - 1$$

Least solution:

$$Z = \begin{cases} 1 & \text{if } 0 \leq q \leq 0.5, \\ \frac{1}{q} - 1 & \text{otherwise.} \end{cases}$$

Training: Rule Extraction

Word-aligned Training Pair

⟨30 duonianlai de youhao hezuo, over 30 years of friendly cooperation⟩

	friendly	cooperation	over	the	last	30	years
30						■	
duonianlai			■	■	■		■
de							
youhao	■						
hezuo		■					

Figure 2 of (Chiang 2007)

Training: Rule Extraction

Initial Rules

- X → ⟨30 duonianlai de youhao hezuo, over 30 years of friendly cooperation⟩
- X → ⟨duonianlai, over the last 30 years⟩
- X → ⟨youhao hezuo, friendly cooperation⟩
- X → ⟨30, 30⟩
- X → ⟨youhao, friendly⟩
- X → ⟨hezuo, cooperation⟩

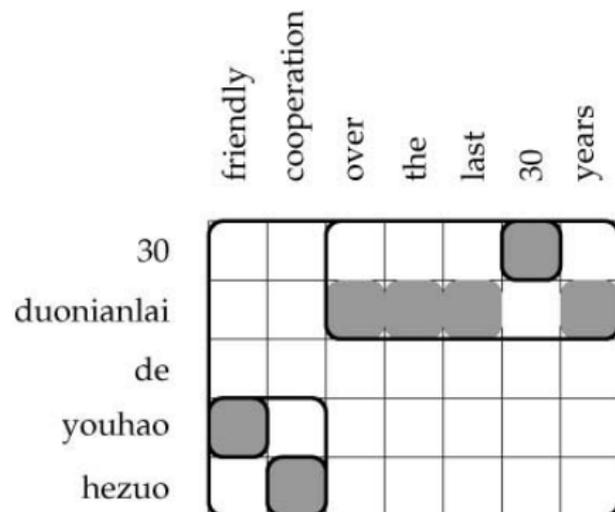


Figure 2 of (Chiang 2007)

Training: Rule Extraction

Example Rule

$X \rightarrow \langle X_{[1]} \text{ duonianlai de } X_{[2]}, X_{[2]} \text{ over the last years } X_{[1]} \rangle$

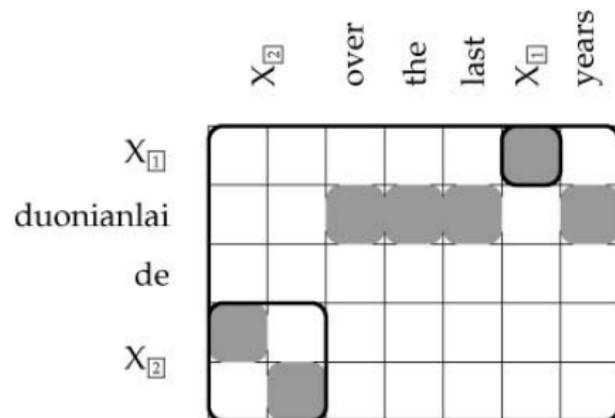


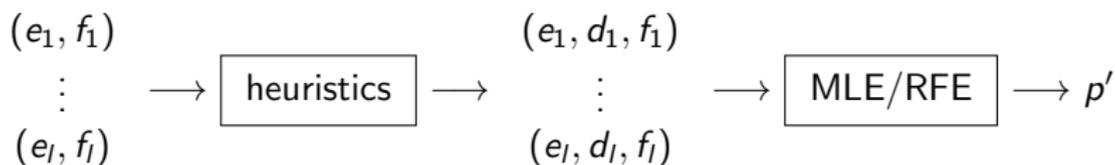
Figure 2 of (Chiang 2007)

Training: Maximum-Likelihood Estimation

- ▶ training data: sequence $\mathcal{D} = (e_1, f_1), \dots, (e_l, f_l)$
- ▶ model: $\mathcal{M}_G = \{p \in \mathcal{M}(E, F) \mid \exists p' : p(e, f) = \sum_d p'(e, d, f)\}$
- ▶ likelihood: $P(\mathcal{D} \mid p') = \prod_j \sum_d p'(e_j, d, f_j)$
- ▶ use EM algorithm

Training: Reality

Chiang (2007), DeNeeffe and Knight (2009):



Outline

Statistical Machine Translation

Weighted Grammars as a Feature

Statistical Machine Translation with Weighted Grammars

Feature Selection

Hidden variable: derivation of (e, f) in G

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
synchronous cfg

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
synchronous cfg
- ▶ $h_3(e, f) = \log \sum_a p(f, a | e)$
alignment-based model

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
synchronous cfg
- ▶ $h_3(e, f) = \log \sum_a p(f, a | e)$
alignment-based model
- ▶ $h_4(e, f) = -|e|$
length preference

Feature Selection

Hidden variable: derivation of (e, f) in G



Chiang (2007):

- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
synchronous cfg
- ▶ $h_3(e, f) = \log \sum_a p(f, a | e)$
alignment-based model
- ▶ $h_4(e, f) = -|e|$
length preference

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
synchronous cfg
- ▶ $h_3(e, f) = \log \sum_a p(f, a | e)$
alignment-based model
- ▶ $h_4(e, f) = -|e|$
length preference



Chiang (2007):

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$
- ▶ $h_1(e, d, f) = \log p(e)$

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
synchronous cfg
- ▶ $h_3(e, f) = \log \sum_a p(f, a | e)$
alignment-based model
- ▶ $h_4(e, f) = -|e|$
length preference



Chiang (2007):

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$
- ▶ $h_1(e, d, f) = \log p(e)$
- ▶ $h_2(e, d, f) = \log p_G(e, d | f)$

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
synchronous cfg
- ▶ $h_3(e, f) = \log \sum_a p(f, a | e)$
alignment-based model
- ▶ $h_4(e, f) = -|e|$
length preference



Chiang (2007):

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$
- ▶ $h_1(e, d, f) = \log p(e)$
- ▶ $h_2(e, d, f) = \log p_G(e, d | f)$
- ▶ $h_3(e, d, f) = \log p_G(f, d | e)$

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
 n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
 synchronous cfg
- ▶ $h_3(e, f) = \log \sum_a p(f, a | e)$
 alignment-based model
- ▶ $h_4(e, f) = -|e|$
 length preference



Chiang (2007):

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$
- ▶ $h_1(e, d, f) = \log p(e)$
- ▶ $h_2(e, d, f) = \log p_G(e, d | f)$
- ▶ $h_3(e, d, f) = \log p_G(f, d | e)$
- ▶ $h_4(e, d, f) = -|e|$

Feature Selection

Hidden variable: derivation of (e, f) in G



Chiang (2007):

- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
synchronous cfg
- ▶ $h_3(e, f) = \log \sum_a p(f, a | e)$
alignment-based model
- ▶ $h_4(e, f) = -|e|$
length preference

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$
- ▶ $h_1(e, d, f) = \log p(e)$
- ▶ $h_2(e, d, f) = \log p_G(e, d | f)$
- ▶ $h_3(e, d, f) = \log p_G(f, d | e)$
- ▶ $h_4(e, d, f) = -|e|$
- ▶ $h_5(e, d, f) = \text{lexical weights}$

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
synchronous cfg
- ▶ $h_3(e, f) = \log \sum_a p(f, a | e)$
alignment-based model
- ▶ $h_4(e, f) = -|e|$
length preference



Chiang (2007):

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$
- ▶ $h_1(e, d, f) = \log p(e)$
- ▶ $h_2(e, d, f) = \log p_G(e, d | f)$
- ▶ $h_3(e, d, f) = \log p_G(f, d | e)$
- ▶ $h_4(e, d, f) = -|e|$
- ▶ $h_5(e, d, f) = \text{lexical weights}$
- ▶ $h_6(e, d, f) = -\#_{\text{ex}} d$

Feature Selection

Hidden variable: derivation of (e, f) in G



- ▶ log-linear model $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶ $h_1(e, f) = \log p(e)$
n-gram language model
- ▶ $h_2(e, f) = \log \sum_d p_G(e, d, f)$
synchronous cfg
- ▶ $h_3(e, f) = \log \sum_a p(f, a | e)$
alignment-based model
- ▶ $h_4(e, f) = -|e|$
length preference



Chiang (2007):

- ▶ $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$
- ▶ $h_1(e, d, f) = \log p(e)$
- ▶ $h_2(e, d, f) = \log p_G(e, d | f)$
- ▶ $h_3(e, d, f) = \log p_G(f, d | e)$
- ▶ $h_4(e, d, f) = -|e|$
- ▶ $h_5(e, d, f) = \text{lexical weights}$
- ▶ $h_6(e, d, f) = -\#_{\text{ex}} d$
- ▶ $h_7(e, d, f) = -\#_{\text{gl}} d$

Decoding

$$\hat{e} = \operatorname{argmax}_e p(e | f)$$

Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f)\end{aligned}$$

Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)}\end{aligned}$$

Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)} \\ &= \operatorname{argmax}_e \sum_d \exp \sum_i \lambda_i \cdot h_i(e, d, f)\end{aligned}$$

Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)} \\ &= \operatorname{argmax}_e \sum_d \exp \sum_i \lambda_i \cdot h_i(e, d, f)\end{aligned}$$

NP-hard! (Casacuberta and Higuera 2000;
Li, Eisner, and Khudanpur 2009)

Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)} \\ &= \operatorname{argmax}_e \sum_d \exp \sum_i \lambda_i \cdot h_i(e, d, f)\end{aligned}$$

NP-hard! (Casacuberta and Higuera 2000;
Li, Eisner, and Khudanpur 2009)

$$\approx \operatorname{argmax}_e \sum_{d \in D_{\text{fin}}} \exp \sum_i \lambda_i \cdot h_i(e, d, f)$$

Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)} \\ &= \operatorname{argmax}_e \sum_d \exp \sum_i \lambda_i \cdot h_i(e, d, f)\end{aligned}$$

NP-hard! (Casacuberta and Higuera 2000;
Li, Eisner, and Khudanpur 2009)

$$\approx \operatorname{argmax}_e \sum_{d \in D_{\text{fin}}} \exp \sum_i \lambda_i \cdot h_i(e, d, f)$$

where $D_{\text{fin}} = \operatorname{argmax}_d^{(n)} \exp \sum_i \lambda_i \cdot h_i(e(d), d, f)$

Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)} \\ &= \operatorname{argmax}_e \sum_d \exp \sum_i \lambda_i \cdot h_i(e, d, f)\end{aligned}$$

NP-hard! (Casacuberta and Higuera 2000;
Li, Eisner, and Khudanpur 2009)

$$\approx \operatorname{argmax}_e \sum_{d \in D_{\text{fin}}} \exp \sum_i \lambda_i \cdot h_i(e, d, f)$$

where $D_{\text{fin}} = \operatorname{argmax}_d^{(n)} \exp \sum_i \lambda_i \cdot h_i(e(d), d, f)$

Cube Pruning (Chiang 2007)

- Banerjee, Satanjeev and Alon Lavie (2005). "METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments". In: [Proceedings of the ACL 2005 Workshop on Intrinsic and Extrinsic Evaluation Measures for MT and/or Summarization](#).
- Berger, Adam L., Vincent J. Della Pietra, and Stephen A. Della Pietra (1996). "A maximum entropy approach to natural language processing". In: [Comp. Ling.](#) 22 (1), pp. 39–71. ISSN: 0891-2017. URL: <http://portal.acm.org/citation.cfm?id=234285.234289>.
- Brown, Peter F. et al. (1993). "The mathematics of statistical machine translation: parameter estimation". In: [Comp. Ling.](#) 19 (2), pp. 263–311. ISSN: 0891-2017. URL: <http://portal.acm.org/citation.cfm?id=972470.972474>.
- Casacuberta, Francisco and Colin de la Higuera (2000). "Computational Complexity of Problems on Probabilistic Grammars and Transducers". In: [LNCS](#).
- Casella, George and Roger L. Berger (2002). [Statistical Inference](#). Duxbury Advances Series. Duxbury.
- Chiang, David (2007). "Hierarchical Phrase-Based Translation". In: [Comp. Ling.](#) 33.2, pp. 201–228. ISSN: 0891-2017. DOI: <http://dx.doi.org/10.1162/coli.2007.33.2.201>.
- DeNeefe, Steve and Kevin Knight (2009). "Synchronous tree adjoining machine translation". In: [EMNLP '09: Proceedings of the 2009](#)

Conference on Empirical Methods in Natural Language Processing. Singapore: Association for Computational Linguistics, pp. 727–736. ISBN: 978-1-932432-62-6.

- Li, Zhifei, Jason Eisner, and Sanjeev Khudanpur (2009). “Variational decoding for statistical machine translation”. In: [ACL-IJCNLP '09: Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2](#). Suntec, Singapore: Association for Computational Linguistics, pp. 593–601. ISBN: 978-1-932432-46-6.
- Och, Franz Josef (2003). “Minimum Error Rate Training in Statistical Machine Translation”. In: [ACL](#), pp. 160–167.
- Och, Franz Josef and Hermann Ney (2002). “Discriminative Training and Maximum Entropy Models for Statistical Machine Translation”. In: [ACL](#), pp. 295–302.
- Papineni, Kishore et al. (2002). “BLEU: a method for automatic evaluation of machine translation”. In: [ACL '02: Proceedings of the 40th Annual Meeting on Association for Computational Linguistics](#). Philadelphia, Pennsylvania: Association for Computational Linguistics, pp. 311–318. DOI: <http://dx.doi.org/10.3115/1073083.1073135>.
- Smith, David A. and Jason Eisner (2006). “Minimum-Risk Annealing for Training Log-Linear Models”. In: [Proceedings of the International Conference on Computational Linguistics and the Association for](#)

Computational Linguistics (COLING-ACL), Companion Volume.

Sydney, pp. 787–794. URL:

<http://cs.jhu.edu/~jason/papers/#acl06-risk>.

Snover, Matthew et al. (2006). “A Study of Translation Edit Rate with Targeted Human Annotation”. In: [Proceedings of Association for Machine Translation in the Americas](#).