Domain Adaptation in Statistical MT

Barry Haddow

University of Edinburgh

3 October, 2013

SMT at Edinburgh: EU Projects

AC:DEPT







MOSES CORE

SMT at Edinburgh: People

Faculty





Researchers



Philipp Koehn

Lexi Birch

Nadir Durrani Matthias Huck

Hieu Hoang

Students





Christian Buck Maria Nadejde



Miles Osborne









Liane Guillou



Phil Williams





Bonnie Webber Ulrich Germann Annie Louis

Saint-Amand









Eva Hasler

Dominikus Wetzel 3 October, 2013

Barry Haddow (University of Edinburgh)

Domain Adaptation in Statistical MT

3 / 42

The minutes of yesterday's sitting have been distributed. Are there any comments? The next item is the statements by the Council and the Commission on the situation in Palestine. The minutes of yesterday's sitting lake been distributed. Are there any comments? The next item is the statements by the Council and the Commission on the situation in Palestine a Re: Norton Ghost 15 Install errors and "error in Main configeration. nope was during first back-up but I can't even start Ghost process without it crashing with taht eror on startup. Re: Norton Ghost 15 Install errors and "error in Main configeration. nope was during first back process without it crashing with taht eror on startup. Yetta, any chicken left from last night? I'il make you a sandwich. Make it to go. Is Shayna coming or what? Come on, give it to me. Yetta, any chicken left from last night? l'il make you a sandwich. Make it to go. Is Shayna coming or what? Come on, give it to me.

- Encompasses notions of genre and topic.
- Practical definition as subcorpus
- But
 - Domains such as news cover many topics
 - Some domains more different than others
- Main issue is difference between training and test data

The SMT Training Pipeline



The Standard (Phrase-Based) SMT Model

Best translation given by:

$$e^* = rgmax_e \sum_i \lambda_i h_i(e, f)$$

Typical features:

- log of forward and reverse translation probabilities
- forward and reverse lexical scores
- log of language model probability
- phrase and word penalty
- reordering model scores

unadapted Europarl \approx 2 million sentences adapted Europarl plus news-commentary (\approx 140k sentences)

test news commentary

What goes wrong? (I)

. . .

Source Als gelernter Arabist , der tief in die arabische und mohammedanische Kultur eingedrungen war ...

Unadapted As a skilled Arabist , deeply in the Arab and mohammedanische culture ...

Adapted As a skilled Arabist steeped in the Arab and Muslim culture ...

Reference As an Arabist by training , immersed in Arab and $\underline{\mathrm{Muslim}}$ culture

What goes wrong? (I)

Source Als gelernter Arabist , der tief in die arabische und mohammedanische Kultur eingedrungen war ...

Unadapted As a skilled Arabist , deeply in the Arab and mohammedanische culture ...

Adapted As a skilled Arabist steeped in the Arab and Muslim culture ...

Reference As an Arabist by training , immersed in Arab and $\underline{\mathrm{Muslim}}$ culture

Occurrences of mohammedanische:

europarl: 0 news-commentary: 13

. . .

Seen Error

Source die indische Zentralbank schwimmt im Geld Unadapted the Indian Central Bank in the money they swim Adapted India's central bank is awash in money Reference India's central bank is rolling in cash Source die indische Zentralbank schwimmt im Geld Unadapted the Indian Central Bank in the money they swim Adapted India's central bank is awash in money Reference India's central bank is rolling in cash

Translations of schwimmt

Unadapted		Adapted	
swim	2/5	swim	2/11
floats	1/5	is awash	2/11
then something	1/5	awash	1/11
then	1/5	rows	1/11

Sense Error

What goes wrong? (III)

Source ... suchte Japan auf breiter Front nach Hilfsmitteln und Innovationen ...

Unadapted \ldots Japan on a broad front sought after aids and innovation \ldots

Adapted ... Japan on a broad front sought after tools and innovations ...

Reference ... Japan searched broadly for tools and innovations ...

What goes wrong? (III)

Source ... suchte Japan auf breiter Front nach Hilfsmitteln und Innovationen ...

Unadapted ... Japan on a broad front sought after aids and innovation ... Adapted ... Japan on a broad front sought after tools and innovations ...

Reference \ldots Japan searched broadly for tools and innovations \ldots

Translations of Hilfsmitteln

Unada	pted	Adapted	
aid	8/38	tools	10/46
tools	6/38	aid	8/46
aids	4/38	aids	4/46
devices	3/38	resources	3/46

Score Error











Experimental Setup

out-of-domain ep European parliament in-domain nc News commentary st Subtitles

- Language and reordering model built from all data
- Tune and test on in-domain
- Average across 8 language pairs

Supplementing In-Domain Data

30 +1.0 +0.5 +0.5 +0.3 +0.6 +0.3 -0.2 20 25 20 15 Bleu Bleu 15 10 10 5 5 0 0 g nc+epA nc+epS nc+epE ep+nc st st+epA st+epS st+epE

News Commentary

Subtitles

Barry Haddow (University of Edinburgh)

Domain Adaptation in Statistical MT

3 October, 2013 16 / 42

ep+st

+0.2

Supplementing Out-of-Domain Data



News Commentary

Subtitles



Barry Haddow (University of Edinburgh)

Domain Adaptation in Statistical MT

3 October, 2013 17 / 42

Effect of Adding Data - Conclusions



- Adding out-of-domain appears to help with OOVs?
 - What about other frequencies?
- Measure source-word precision by tracking word translation
- Calculate average precision, binned by log frequency (in training).

Precision vs Frequency (NC)



Precision vs Frequency (ST)



Precision vs Frequency - Conclusions



Methods for Domain Adaptation in SMT





• Instead of using all training data ...



- Instead of using all training data ...
- Use only selected data
- But how to select?



- Instead of using all training data ...
- Use only selected data
- But how to select?
 - Modified Moore-Lewis
 - Select similar to in-domain ...
 - ...but different from out
 - Use LM perplexity for similarity



• Instance selection is a 1-0 weighting



- Instance selection is a 1-0 weighting
- Can we improve by a allowing variable weights?
 - Use MML scores for weighting
 - .. really should learn weights

WMT13 Experiments

Parallel Training Data



Test data: Extracted from online news.

Instance Selection and Weighting - WMT13 Experiments





Mixture Model





Which models to mix?



Log-Linear $\dots w_{T_1} \log(p_1(e|f)) + w_{T_2} \log(p_2(e|f)) \dots$

- Optimise with MERT etc.
- Problems with zeros





- More natural
- Separate optimisation

Linear Mixture Models - WMT13 Experiments



- Linear mixture better than log linear
 - $\rightarrow~$ But perplexity is indirect and ill-defined objective
- What if we could optimise directly for translation performance?
 e.g. BLEU
- This can be done with Pairwise Ranked Optimisation (PRO)

PRO: Pairwise Ranked Optimisation

- Batch tuning algorithm optimises standard linear model
- Sample pairs of hypotheses from *n*-best lists

$$S = \{(e_1^1, e_2^1), \dots (e_1^n, e_2^n)\}$$

• Feature weights optimised by:

$$\begin{split} \mathbf{w}^* &= \arg \max_{\mathbf{w}} \quad \sum_{i=1}^n \log \left(\sigma \left(y_i \cdot (score_{\mathbf{w}}(e_1^i) - score_{\mathbf{w}}(e_2^i)) \right) \right) \\ & \text{where} \qquad \sigma(x) = 1/(1 + \exp(-x)) \\ & \text{and} \qquad y_i = \operatorname{sgn} \left(\operatorname{bleu}(e_1^i) - \operatorname{bleu}(e_2^i) \right) \end{split}$$

PRO for Mixture Model Training

Standard linear model

$$\mathsf{score}_{\mathsf{w}}(e) = \mathsf{h}(e) \cdot \mathsf{w}$$

Interpolated TM makes score() function of mixture weights

$$\sum_{j=1}^{m} \left(w^j \cdot \log \left(\lambda^j p_A^j + (1 - \lambda^j) p_B^j \right) \right) + \sum_{j=m+1}^{n} w^j h^j$$

• So just optimise PRO objective for both w and λ !

Experimental Setup

out-of-domain	ер	European parliament
in-domain	nc	News commentary
	st	Subtitles

- Language and reordering model on all data
- Tune and test on in-domain
- Average across 8 language pairs
- Baseline is concatenation

PRO for Mixture Models - Results



Barry Haddow (University of Edinburgh)

Domain Adaptation in Statistical MT

3 October, 2013 35 / 42

Feature Engineering for Domain Adaptation

- Can now tune models with many features
- Can this help with domain adaptation?

Feature Engineering for Domain Adaptation

- Can now tune models with many features
- Can this help with domain adaptation?

	Standard	Word	Word-Topic
matter	p _{dir} =0.2,		
\updownarrow	p _{inv} =0.4,	wp_matter_important=1	wt_matter_important_T1 ${=}1$
important	l _{dir} =0.5,		wt_matter_important_T6 ${=}1$
	$l_{inv} = 0.1$		
	p _{dir} =0.1,		
matter	p _{inv} =0.6,		wt_matter_matiere_T2 ${=}1$
\uparrow	<i>l</i> _{dir} =0.4,	wp_matter_matiere $=1$	wt_matter_matiere_T7 ${=}1$
matière	l _{inv} =0.2		

Feature Engineering: Results



Experiments using IWSLT (TED) Data

Training Data

Madame la Présidente, c'est une motion de procédure. Vous avez probablement appris par la presse et par la télévision que plusieurs attentats à la bombe et crimes ont été perpétrés au Sri Lanka.

Support Forums

Tu ne retrouve pas ton compte Norton, je te conseille de joindre le Support des produits Norton et ils seront à meme via différents moyens, de le retrouver pour toi.

Training Data

Madame la Présidente, c'est une motion de procédure. Vous avez probablement appris par la presse et par la télévision que plusieurs attentats à la bombe et crimes ont été perpétrés au Sri Lanka.

Support Forums

Tu ne retrouve pas ton compte Norton, je te conseille de joindre le Support des produits Norton et ils seront à meme via différents moyens, de le retrouver pour toi.

How can we deal with register differences?

Informal vs. Formal French

Tu l'as dit toi-même Est-ce que tu as des ... ↔ Vous l'avez dit vous-même
 ↔ Avez-vous des ...

• Two Approaches:

- \rightarrow Transform the training data
- ightarrow Transform the test data
- Both methods are effective for *tu/vous*
 - Transform test better, but effect is additive
- est-ce que transform not effective too much variety

- Heuristic, multi-stage pipeline makes DA difficult
- OOVs biggest problem, also score errors
- Variety of techniques:
 - Looked mainly at filtering and model interpolation
 - Really only tackling score errors
- Often no clear winner ... use many languages and data sets

- Better analysis of why things work
- Move away from data-set as domain
- Deal better with informal text preprocessing
- Improve handling of OOVs
- Make more use of non-parallel data
- Improved pipeline could make DA easier to analyse

Thank You. Questions?

Collaborators

Pierrette Bouillon, Nadir Durrani, Eva Hasler, Kenneth Heafield, Philipp Koehn, Manny Rayner