ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

# Tuning SMT Systems on the Training Set

Chris Dyer, Patrick Simianer, Stefan Riezler, Phil Blunsom,
Eva Hasler

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

Goal: Discriminative training using **sparse features** on the **full training set**

Goal: Discriminative training using **sparse features** on the **full training set**

Approach: Picky-picky / elitist learning:

Goal: Discriminative training using **sparse features** on the **full training set**

Approach: Picky-picky / elitist learning:
- Stochastic learning with **true random selection of examples**.

# Tuning SMT Systems on the Training Set

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

Goal: Discriminative training using **sparse features** on the **full training set**

Approach: Picky-picky / elitist learning:

- Stochastic learning with **true random selection of examples**.
- **Feature selection** according to various regularization criteria.

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

Goal: Discriminative training using **sparse features** on the **full training set**

Approach: Picky-picky / elitist learning:

- Stochastic learning with **true random selection of examples**.
- **Feature selection** according to various regularization criteria.
- **Leave-one-out estimation**: Leave out sentence/shard currently being trained on when extracting rules/features in training.

- cdec decoder (https://github.com/redpony/cdec)

- cdec decoder (https://github.com/redpony/cdec)
- Hiero SCFG grammars

# SMT Framework + Data

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- cdec decoder (https://github.com/redpony/cdec)
- Hiero SCFG grammars
- WMT11 news-commentary corpus

# SMT Framework + Data

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- cdec decoder (https://github.com/redpony/cdec)
- Hiero SCFG grammars
- WMT11 news-commentary corpus
  - 132,755 parallel sentences

# SMT Framework + Data

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- cdec decoder (https://github.com/redpony/cdec)
- Hiero SCFG grammars
- WMT11 news-commentary corpus
    - 132,755 parallel sentences
    - German-to-English

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

**Algorithm** extended ranking voted perceptron: training

$D = \{D^1, ..., D^M\}$: Development set

$C^m = \{c_1^m, ..., c_N^m\}$: the original $N$-best list of $D^m$

$c_n^m$: $n$-th candidate in $C^m$

$X^m = \{x_1^m, ...x_N^m\}$: (reordered) $N$-best list of $D^m$

$x_i^m$: $i$-th candidate in the (reordered) $N$-best list $X^m$

$Ranking(W, C^m)$: returns $N$-best list of $C^m$ reordered
    based on the score, $s_n^m = <W, \phi(c_n^m)>$

$\phi(x_n^m)$: the feature vector of $x_n^m$

$W$: weight vector

$V = \{V_1, ...V_T\}$: set of weight vectors

$T$: Number of pre-defined iteration

1: **For** $t = 1, ..., T$
2:   **For** $m = 1, ..., M$ ;; for each sample in dev-set
3:     $X^m \leftarrow Ranking(W, C^m)$
4:     **For** $i = 1, ..., |X^m|$
5:       **For** $j = i + 1, ..., |X^m|$
6:         **If** $(BLEU(x_j^m) > BLEU(x_i^m)$
7:             $\& \ WER(x_j^m) <= WER(x_i^m))$
8:           $s = (BLEU(x_j^m) - BLEU(x_i^m))$
9:           $W = W + s * (\phi(x_j^m) - \phi(x_i^m))$
10:        **End_If**
11:      **End_For**
12:    **End_For**
13:    $V_t = W$
14:  **End_For**
15: **End_For**
16: **Return** $V$

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Random sampling of pairs from full chart for pairwise ranking:

- Random sampling of pairs from full chart for pairwise ranking:
  - First sample translations according to their model score.

- Random sampling of pairs from full chart for pairwise ranking:
    - First sample translations according to their model score.
    - Then sample pairs.

- Random sampling of pairs from full chart for pairwise ranking:
    - First sample translations according to their model score.
    - Then sample pairs.
- Sampling will diminish problem of learning to discriminate translations that are too close (in terms of sentence-wise approx. BLEU) to each other.

# Constraint Selection = Sampling of Pairs

- Random sampling of pairs from full chart for pairwise ranking:
    - First sample translations according to their model score.
    - Then sample pairs.
- Sampling will diminish problem of learning to discriminate translations that are too close (in terms of sentence-wise approx. BLEU) to each other.
- Sampling will also speed up learning.

# Constraint Selection = Sampling of Pairs

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Random sampling of pairs from full chart for pairwise ranking:
    - First sample translations according to their model score.
    - Then sample pairs.
- Sampling will diminish problem of learning to discriminate translations that are too close (in terms of sentence-wise approx. BLEU) to each other.
- Sampling will also speed up learning.
- Lots of variations on sampling possible ...

# Candidate Features

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Efficient computation of features from local rule context:

# Candidate Features

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Efficient computation of features from local rule context:
  - Hiero SCFG rule identifier

# Candidate Features

- Efficient computation of features from local rule context:
  - Hiero SCFG rule identifier
  - target n-grams within rule

# Candidate Features

- Efficient computation of features from local rule context:
    - Hiero SCFG rule identifier
    - target n-grams within rule
    - target n-gram with gaps (X) within rule

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

# Candidate Features

- Efficient computation of features from local rule context:
  - Hiero SCFG rule identifier
  - target n-grams within rule
  - target n-gram with gaps (X) within rule
  - binned rule counts in full training set

# Candidate Features

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Efficient computation of features from local rule context:
  - Hiero SCFG rule identifier
  - target n-grams within rule
  - target n-gram with gaps (X) within rule
  - binned rule counts in full training set
  - rule length features

# Candidate Features

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Efficient computation of features from local rule context:
  - Hiero SCFG rule identifier
  - target n-grams within rule
  - target n-gram with gaps (X) within rule
  - binned rule counts in full training set
  - rule length features
  - rule shape features

# Candidate Features

- Efficient computation of features from local rule context:
  - Hiero SCFG rule identifier
  - target n-grams within rule
  - target n-gram with gaps (X) within rule
  - binned rule counts in full training set
  - rule length features
  - rule shape features
  - word alignments in rules

# Candidate Features

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Efficient computation of features from local rule context:
  - Hiero SCFG rule identifier
  - target n-grams within rule
  - target n-gram with gaps (X) within rule
  - binned rule counts in full training set
  - rule length features
  - rule shape features
  - word alignments in rules
- ... and many more!

# Feature Selection

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- $\ell_1/\ell_2$-regularization

# Feature Selection

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- $\ell_1/\ell_2$-regularization
  - Compute $\ell_2$-norm of column vectors ($=$ vector of examples/shards for each of $n$ features), then $\ell_1$-norm of resulting $n$-dimensional vector.

# Feature Selection

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- $\ell_1/\ell_2$-regularization
  - Compute $\ell_2$-norm of column vectors ($=$ vector of examples/shards for each of $n$ features), then $\ell_1$-norm of resulting $n$-dimensional vector.

$$\mathbf{W_a} : \begin{bmatrix} 4 & 0 & 0 & 3 \\ 0 & 4 & 3 & 0 \end{bmatrix} \quad \mathbf{W_b} : \begin{bmatrix} 4 & 3 & 0 & 0 \\ 0 & 4 & 3 & 0 \end{bmatrix}$$
$$\quad 4 \quad 4 \quad 3 \quad 3 \rightarrow 14 \qquad 4 \quad 5 \quad 3 \quad 0 \rightarrow 12$$

# Feature Selection

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- $\ell_1/\ell_2$-regularization
  - Compute $\ell_2$-norm of column vectors ($=$ vector of examples/shards for each of $n$ features), then $\ell_1$-norm of resulting $n$-dimensional vector.

$$\mathbf{W}_a : \begin{bmatrix} 4 & 0 & 0 & 3 \\ 0 & 4 & 3 & 0 \end{bmatrix} \quad \mathbf{W}_b : \begin{bmatrix} 4 & 3 & 0 & 0 \\ 0 & 4 & 3 & 0 \end{bmatrix}$$
$$4 \quad 4 \quad 3 \quad 3 \to 14 \qquad 4 \quad 5 \quad 3 \quad 0 \to 12$$

- Effect is to choose small subset of features that are useful across all examples/shards

# Feature Selection, done properly

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Incremental gradient-based selection of column vectors (Obozinski, Taskar, Jordan: Joint covariant selection and joint subspace selection for multiple classification problems. Stat Comput (2010))

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

# Feature Selection, done properly

- Incremental gradient-based selection of column vectors (Obozinski, Taskar, Jordan: Joint covariant selection and joint subspace selection for multiple classification problems. Stat Comput (2010))

---

**Algorithm 1** Approximate block-Lasso path

Given $\epsilon$ and $\xi$,
**while** $\lambda^t > \lambda_{\min}$ **do**

Set $j* = \operatorname{argmax}_j \|\nabla_{w_j} J(W^t)\|$

Update $w_{j*}^{(t+1)} = w_{j*}^{(t)} - \epsilon u^t$ with $u^t = \frac{\nabla_{w_{j*}} J}{\|\nabla_{w_{j*}} J\|}$

$\lambda^{t+1} = \min\left(\lambda^t, \frac{J(W^t) - J(W^{t+1})}{\epsilon}\right)$

Add $j*$ to the active set

Enforce (4) for covariates in the active set with $\xi_0 = \xi$.

**end while**

---

- Combine feature selection with averaging:

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Combine feature selection with averaging:
  - Keep only those features with large enough $\ell_2$-norm computed over examples/shards.

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Combine feature selection with averaging:
  - Keep only those features with large enough $\ell_2$-norm computed over examples/shards.
  - Then average feature values over examples/shards.

- First full training run finished!

- First full training run finished!
    - 150k parallel sentences from news commentary data,
    German-to-English

- First full training run finished!
  - 150k parallel sentences from news commentary data,
    German-to-English
  - pairwise ranking perceptron

- First full training run finished!
  - 150k parallel sentences from news commentary data, German-to-English
  - pairwise ranking perceptron
  - sample 100 translations from chart, use all $100*(99)/2$ pairs

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- First full training run finished!
    - 150k parallel sentences from news commentary data, German-to-English
    - pairwise ranking perceptron
    - sample 100 translations from chart, use all 100*(99)/2 pairs
    - OR: use n-best list
    - sparse rule-id features AND/OR dense features

- First full training run finished!
  - 150k parallel sentences from news commentary data, German-to-English
  - pairwise ranking perceptron
  - sample 100 translations from chart, use all 100*(99)/2 pairs
  - OR: use n-best list
  - sparse rule-id features AND/OR dense features
  - 200 shards (25 machines with 8 cores)

- Still a lot of bugs due to integration of code from different sources

- Still a lot of bugs due to integration of code from different sources
- Infrastructure is working

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Still a lot of bugs due to integration of code from different sources
- Infrastructure is working
- Experiments still running

# Results

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Still a lot of bugs due to integration of code from different sources
- Infrastructure is working
- Experiments still running
- Sensible things happening:
  - Best rule $X \rightarrow X_1$ , $\mathrm{dass}$ $X_2$, $X_1$ $\mathrm{that}$ $X_2$
  - Bad rule $X \rightarrow X_1$ $\mathrm{oder}$ $X_2$, $X_1$ $\mathrm{and}$ $X_2$

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

## Results

- Still a lot of bugs due to integration of code from different sources
- Infrastructure is working
- Experiments still running
- Sensible things happening:
    - Best rule $X \to X_1$ , dass $X_2$, $X_1$ that $X_2$
    - Bad rule $X \to X_1$ oder $X_2$, $X_1$ and $X_2$
- At the moment still trailing behind MERT ...

# Results

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

- Still a lot of bugs due to integration of code from different sources
- Infrastructure is working
- Experiments still running
- Sensible things happening:
  - Best rule $X \rightarrow X_1$ , dass $X_2$, $X_1$ that $X_2$
  - Bad rule $X \rightarrow X_1$ oder $X_2$, $X_1$ and $X_2$
- At the moment still trailing behind MERT ...
- We'll catch up!

ToTS

Dyer,
Simianer,
Riezler,
Blunsom,
Hasler

Thanks to organizers for great
opportunity to learn/chat/hobnob!