

SwitchOut: An Efficient Data Augmentation for Neural Machine Translation

Xinyi Wang*, Hieu Pham*, Zihang Dai, Graham Neubig

Carnegie Mellon University



Language
Technologies
Institute

November 2, 2018

*:equal contribution

- Neural models are data hungry, while collecting data is expensive

¹image source:Medium

- Neural models are data hungry, while collecting data is expensive
- Prevalent in computer vision¹

¹image source:Medium

- Neural models are data hungry, while collecting data is expensive
- Prevalent in computer vision¹

- More difficult for natural language
 - | Discrete vocabulary
 - | NMT sensitive to arbitrary noise

¹image source:Medium

Word replacement

Word replacement

- Dictionary [Fadaee et al., 2017]

Word replacement

- Dictionary [Fadaee et al., 2017]
- Word dropout [Sennrich et al., 2016a]

Word replacement

- Dictionary [Fadaee et al., 2017]
- Word dropout [Sennrich et al., 2016a]
- Reward Augmented Maximum Likelihood (RAML) [Norouzi et al., 2016]

Word replacement

- Dictionary [Fadaee et al., 2017]
- Word dropout [Sennrich et al., 2016a]
- Reward Augmented Maximum Likelihood (RAML)
[Norouzi et al., 2016]

! Can we characterize all of the related approaches together?

RAML [Norouzi et al., 2016]

- Motivation: NMT relies on imperfect partial translation at test time, but trained only on gold standard target

RAML [Norouzi et al., 2016]

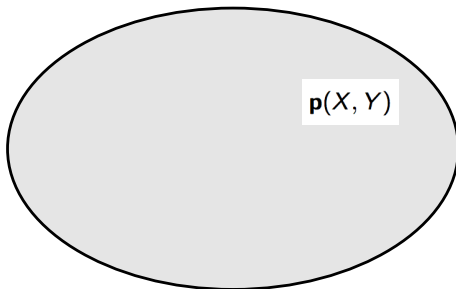
- Motivation: NMT relies on imperfect partial translation at test time, but trained only on gold standard target
- Solution: Sample corrupted target during training

RAML [Norouzi et al., 2016]

- Motivation: NMT relies on imperfect partial translation at test time, but trained only on gold standard target
- Solution: Sample corrupted target during training
- Gold target y , corrupted \tilde{y} , similarity measure r_y

$$q^*(\tilde{y}; y) = \frac{\exp f_{r_y}(\tilde{y}; y) = g}{\sum_{\tilde{y}'} \exp f_{r_y}(\tilde{y}'; y) = g}$$

- **Real** data distribution: $x; y$ $\mathbf{p}(X; Y)$



Formalize Data Augmentation

Real data distribution: $x; y \quad p(X; Y)$

Observed data distribution: $x; y \quad \tilde{p}(X; Y)$

Formalize Data Augmentation

Real data distribution: $x; y \quad p(X; Y)$

Observed data distribution: $x; y \quad \hat{p}(X; Y)$

! Problem: $p(X; Y)$ and $\hat{p}(X; Y)$ might have large discrepancy

Formalize Data Augmentation

Real data distribution: $x; y \quad p(X; Y)$

Observed data distribution: $x; y \quad \hat{p}(X; Y)$

! Problem: $p(X; Y)$ and $\hat{p}(X; Y)$ might have large discrepancy

Data augmentation: $x; y \quad q(X; Y)$

Design a good $q(\mathcal{X}; \Psi)$

q : function of observed $(x; y)$

Design a good $q(\mathcal{X}; \Psi)$

q : function of observed $(x; y)$
How should q approximate p ?

Design a good $q(\mathcal{X}; \Psi)$

q : function of observed $(x; y)$

How should q approximate p ?

- | Diversity: larger support with all valid data pairs $x(y)$
 - ⊆ Entropy $H q(\mathbf{b}; \mathbf{b} | x; y)$ is large

Design a good $q(\mathcal{X}; \Psi)$

q : function of observed $(x; y)$

How should q approximate p ?

- | Diversity: larger support with all valid data pairs $x(y)$
 - Entropy $H q(\mathbf{k}; \Phi; x; y)$ is large
- | Smoothness: probability of similar data pairs are similar
 - q maximizes similarity measure $r_x(x; \mathbf{k}), r_y(y; \Phi)$

Design a good $q(\mathcal{X}; \psi)$

q : function of observed $(x; y)$

How should q approximate p ?

- | Diversity: larger support with all valid data pairs $x(y)$
 - F Entropy $H q(\mathbf{k}; \psi | x; y)$ is large
- | Smoothness: probability of similar data pairs are similar
 - F q maximizes similarity measure $r_x(x; \mathbf{k}), r_y(y; \psi)$

: control effect of diversity; q should maximize

$$J(q) = H q(\mathbf{k}; \psi | x; y) + E_{\mathbf{k}; \psi} q [r_x(x; \mathbf{k}) + r_y(y; \psi)]$$

Mathematically Optimal

$$J(q) = H(q; \theta; \phi; x; y) + E_{\theta; \phi} [q(r_x(x; \theta) + r_y(y; \phi))]$$

Solve for the best q

$$q(\theta; \phi; x; y) = P \frac{\exp(s(\theta; \phi; x; y) = g)}{\int \exp(s(\theta^0; \phi^0; x; y) = g)}$$

Mathematically Optimal

$$J(q) = \int H(q; \theta; x; y) + E_{q; \theta} [r_x(x; \theta) + r_y(y; \theta)]$$

Solve for the best q

$$q(\theta; \theta; x; y) = \frac{\int \exp(s(\theta; \theta; x; y) - g) \exp(s(\theta^0; \theta^0; x; y) - g)}{\int \exp(s(\theta^0; \theta^0; x; y) - g)}$$

Decompose x and y

$$q(\theta; \theta; x; y) = \frac{\int \exp(r_x(\theta; x) - g_x) \exp(r_y(\theta; y) - g_y)}{\int \exp(r_x(\theta^0; x) - g_x) \int \exp(r_y(\theta^0; y) - g_y)}$$

Mathematically Optimal

$$J(q) = \mathbb{E}_{q(\mathbf{k}; \mathbf{p}; x; y)} [r_x(x; \mathbf{k}) + r_y(y; \mathbf{p})]$$

Solve for the best q

$$q(\mathbf{k}; \mathbf{p}; x; y) = \mathbb{P}_{\mathbf{k}^0, \mathbf{p}^0} \frac{\exp s(\mathbf{k}; \mathbf{p}; x; y) = g}{\exp s(\mathbf{k}^0; \mathbf{p}^0; x; y) = g}$$

Decompose x and y

$$q(\mathbf{k}; \mathbf{p}; x; y) = \mathbb{P}_{\mathbf{k}^0} \frac{\exp r_x(\mathbf{k}; x) = xg}{\exp r_x(\mathbf{k}^0; x) = xg} \quad \mathbb{P}_{\mathbf{p}^0} \frac{\exp r_y(\mathbf{p}; y) = yg}{\exp r_y(\mathbf{p}^0; y) = yg}$$

Formulate existing methods

- | Dictionary: jointly on x and y , but deterministic and not diverse
- | Word dropout: only x side with null token
- | RAML: only y side

Formulate SwitchOut

Augment both x and y !

Formulate SwitchOut

Augment both x and y !
Sample for x , y independently

Formulate SwitchOut

Augment both x and y !

Sample for x , y independently

Define $r_x(\mathbf{k}; x)$ and $r_y(\mathbf{p}; y)$

- ▮ Negative Hamming Distance, following RAML

SwitchOut: Sample efficiently

Given a sentence $\theta = \{s_1; s_2; \dots; s_{s_j}\}$

How many words to corrupt?

Assumption: only one token for swapping.

$$P(n) / \exp(-n) =$$

SwitchOut: Sample efficiently

Given a sentence $\theta = \{s_1; s_2; \dots; s_{|s|}\}$

How many words to corrupt?

Assumption: only one token for swapping.

$$P(n) / \exp(-n) =$$

What is the corrupted sentence?

$$P(\text{randomly swap } s_i \text{ by another word}) = \frac{n}{|s|}$$

See Appendix: Efficient batch implementation in PyTorch and TensorFlow

Experiments

Datasets

- | en-vi: IWSLT 2015
- | de-en: IWSLT 2016
- | en-de: WMT 2015

Models

- | Transformer model
- | Word-based, standard preprocessing

Results: RAML and word dropout

src	Method	trg	en-de	de-en	en-vi
N/A		N/A	21.73	29.81	27.97
WordDropout		N/A	20.63	29.97	28.56
SwitchOut		N/A	22.78 ^y	29.94	28.67 ^y
N/A		RAML	22.83	30.66	28.88
WordDropout		RAML	20.69	30.79	28.86
SwitchOut		RAML	23.13 ^y	30.98 ^y	29.09

Results: RAML and word dropout

SwitchOut on source \Rightarrow word dropout

src	Method	trg	en-de	de-en	en-vi
N/A		N/A	21.73	29.81	27.97
	WordDropout	N/A	20.63	29.97	28.56
	SwitchOut	N/A	22.78 ^y	29.94	28.67 ^y
N/A		RAML	22.83	30.66	28.88
WordDropout		RAML	20.69	30.79	28.86
SwitchOut		RAML	23.13 ^y	30.98 ^y	29.09

Results: RAML and word dropout

SwitchOut on source → word dropout

SwitchOut on source and target → RAML

src	Method	trg	en-de	de-en	en-vi
N/A		N/A	21.73	29.81	27.97
	WordDropout	N/A	20.63	29.97	28.56
	SwitchOut	N/A	22.78 ^y	29.94	28.67 ^y
	N/A	RAML	22.83	30.66	28.88
	WordDropout	RAML	20.69	30.79	28.86
	SwitchOut	RAML	23.13 ^y	30.98 ^y	29.09

Where does SwitchOut help?

- More gain for sentences more different from training data

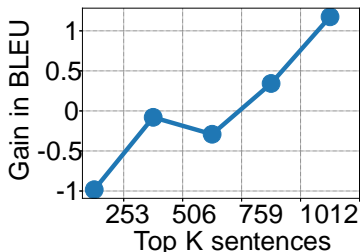
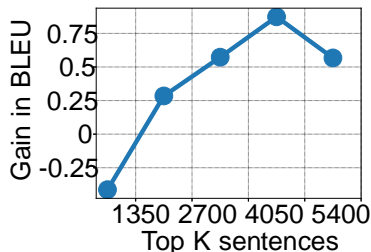


Figure: *Left*: IWSLT 16 de-en. *Right*: IWSLT 15 en-vi.






- SwitchOut sampling is efficient and easy-to-use

- SwitchOut sampling is efficient and easy-to-use
- Work with any NMT architecture

- SwitchOut sampling is efficient and easy-to-use
- Work with any NMT architecture
- Formulation of data augmentation encompasses existing works and inspires future direction

- SwitchOut sampling is efficient and easy-to-use
- Work with any NMT architecture
- Formulation of data augmentation encompasses existing works and inspires future direction

Thanks a lot for listening! Questions?

-  Norouzi et al. (2016) Reward Augmented Maximum Likelihood for Neural Structured Prediction. In NIPS.
-  Sennrich et al. (2016a) Edinburgh neural machine translation systems for wmt 16. In WMT.
-  Sennrich et al. (2016b) Improving neural machine translation models with monolingual data. In ACL.
-  Currey et al. (2017) Copied Monolingual Data Improves Low-Resource Neural Machine Translation. In WMT.
-  Fadaee et al. (2017) Data Augmentation for Low-Resource Neural Machine Translation. In ACL.

- SwitchOut > Back Translation

Method	en-de
Transformer	21.73
+SwitchOut	22.78
+BT	21.82

- SwitchOut > Back Translation
- Switchout + RAML + back translate wins

Method	en-de
Transformer	21.73
+SwitchOut	22.78
+BT	21.82
+BT +RAML	21.53
+BT +SwitchOut	22.93
+BT +RAML +SwitchOut	23.76