

SwitchOut: An Efficient Data Augmentation for Neural Machine Translation

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

Carnegie Mellon University



Language
Technologies
Institute

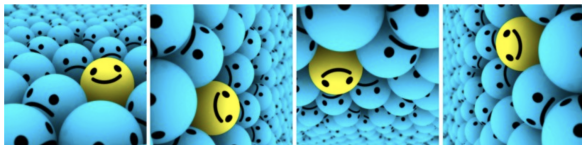
November 2, 2018

*:equal contribution

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

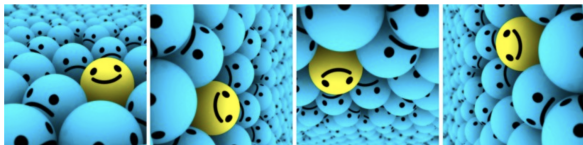
¹image source:Medium

- 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

wie geht es dir ? → How are you ?

Word replacement

wie geht es dir ? → How are you ?

- Dictionary [Fadaee et al., 2017]

wie geht es Tom ? → How is Tom ?

Word replacement

wie geht es dir ? → How are you ?

- Dictionary [Fadaee et al., 2017]

wie geht es Tom ? → How is Tom ?

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

wie geht es NULL ? → How are you ?

Word replacement

wie geht es dir ? \longrightarrow How are you ?

- Dictionary [Fadaee et al., 2017]

wie geht es Tom ? \longrightarrow How is Tom ?

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

wie geht es NULL ? \longrightarrow How are you ?

- Reward Augmented Maximum Likelihood (RAML)
[Norouzi et al., 2016]

wie geht es dir ? \longrightarrow How are hello ?

Word replacement

wie geht es dir ? \longrightarrow How are you ?

- Dictionary [Fadaee et al., 2017]

wie geht es Tom ? \longrightarrow How is Tom ?

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

wie geht es NULL ? \longrightarrow How are you ?


- Reward Augmented Maximum Likelihood (RAML)
[Norouzi et al., 2016]

wie geht es dir ? \longrightarrow How are hello ?

\longrightarrow 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

wie	geht	es	dir	?		How	are	you	?
wie	geht	es	dir	?		How	are	you	?
wie	geht	es	dir	?		How	are	you	?
wie	geht	es	dir	?		How	are	hello	?
wie	geht	es	dir	?		How	yes	you	?
wie	geht	es	dir	?		How	him	her	?

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

wie	geht	es	dir	?		How	are	you	?
wie	geht	es	dir	?		How	are	you	?
wie	geht	es	dir	?		How	are	you	?
wie	geht	es	dir	?		How	are	hello	?
wie	geht	es	dir	?		How	yes	you	?
wie	geht	es	dir	?		How	him	her	?

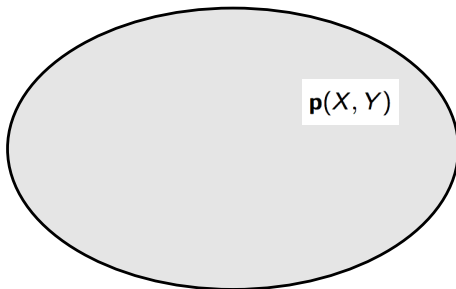
RAML [Norouzi et al., 2016]

- 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 \hat{y} , similarity measure r_y

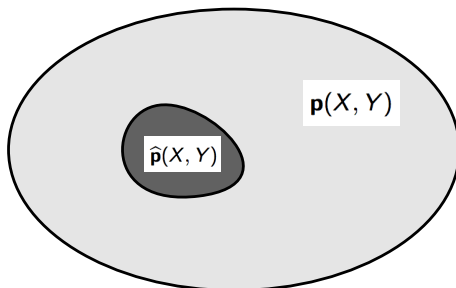
$$\mathbf{q}^*(\hat{y}|y, \tau) = \frac{\exp\{r_y(\hat{y}, y)/\tau\}}{\sum_{\hat{y}'} \exp\{r_y(\hat{y}', y)/\tau\}}$$

wie	geht	es	dir	?		How	are	you	?
wie	geht	es	dir	?		How	are	you	?
wie	geht	es	dir	?		How	are	you	?
wie	geht	es	dir	?		How	are	hello	?
wie	geht	es	dir	?		How	yes	you	?
wie	geht	es	dir	?		How	him	her	?

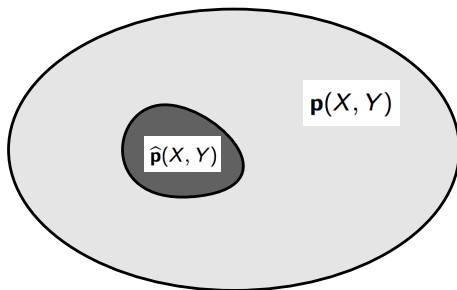
- **Real** data distribution: $x, y \sim \mathbf{p}(X, Y)$



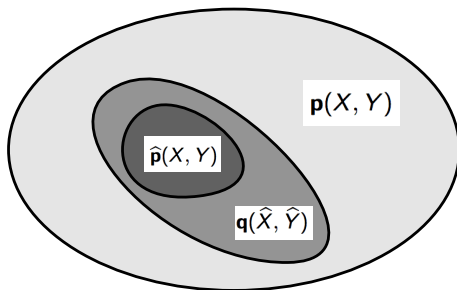
- **Real** data distribution: $x, y \sim \mathbf{p}(X, Y)$
- **Observed** data distribution: $x, y \sim \hat{\mathbf{p}}(X, Y)$



- **Real** data distribution: $x, y \sim \mathbf{p}(X, Y)$
- **Observed** data distribution: $x, y \sim \hat{\mathbf{p}}(X, Y)$
→ Problem: $\mathbf{p}(X, Y)$ and $\hat{\mathbf{p}}(X, Y)$ might have large discrepancy

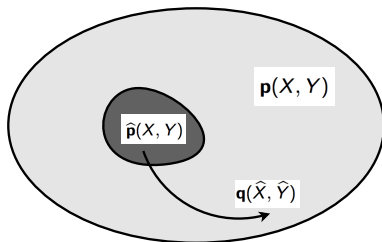


- **Real** data distribution: $x, y \sim \mathbf{p}(X, Y)$
- **Observed** data distribution: $x, y \sim \hat{\mathbf{p}}(X, Y)$
→ Problem: $\mathbf{p}(X, Y)$ and $\hat{\mathbf{p}}(X, Y)$ might have large discrepancy
- Data augmentation: $\hat{x}, \hat{y} \sim \mathbf{q}(\hat{X}, \hat{Y})$



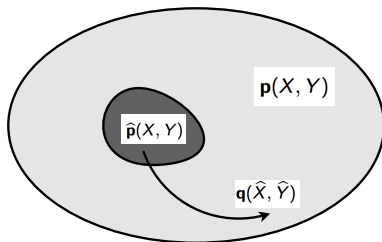
Design a good $\mathbf{q}(\hat{X}, \hat{Y})$

- \mathbf{q} : function of **observed** (x, y)



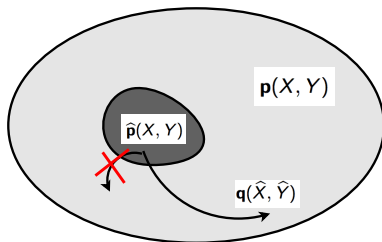
Design a good $q(\hat{X}, \hat{Y})$

- q : function of **observed** (x, y)
- How should q approximate p ?



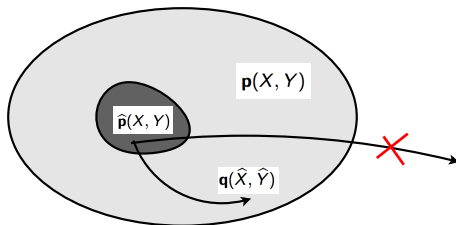
Design a good $q(\hat{X}, \hat{Y})$

- q : function of **observed** (x, y)
- How should q approximate p ?
 - ▶ **Diversity**: larger support with all valid data pairs (x, y)
 - ★ Entropy $\mathbb{H}[q(\hat{x}, \hat{y}|x, y)]$ is large



Design a good $q(\hat{X}, \hat{Y})$

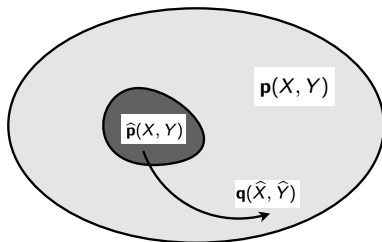
- q : function of **observed** (x, y)
- How should q approximate p ?
 - ▶ **Diversity**: larger support with all valid data pairs (x, y)
 - ★ Entropy $\mathbb{H}[q(\hat{x}, \hat{y}|x, y)]$ is large
 - ▶ **Smoothness**: probability of similar data pairs are similar
 - ★ q maximizes similarity measure $r_x(x, \hat{x}), r_y(y, \hat{y})$



Design a good $\mathbf{q}(\hat{X}, \hat{Y})$

- \mathbf{q} : function of **observed** (x, y)
- How should \mathbf{q} approximate \mathbf{p} ?
 - ▶ **Diversity**: larger support with all valid data pairs (x, y)
 - ★ Entropy $\mathbb{H}[\mathbf{q}(\hat{x}, \hat{y}|x, y)]$ is large
 - ▶ **Smoothness**: probability of similar data pairs are similar
 - ★ \mathbf{q} maximizes similarity measure $r_x(x, \hat{x}), r_y(y, \hat{y})$
- τ : control effect of diversity; \mathbf{q} should maximize

$$J(\mathbf{q}) = \tau \cdot \mathbb{H}[\mathbf{q}(\hat{x}, \hat{y}|x, y)] + \mathbb{E}_{\hat{x}, \hat{y} \sim \mathbf{q}}[r_x(x, \hat{x}) + r_y(y, \hat{y})]$$



$$J(\mathbf{q}) = \tau \cdot \mathbb{H}[\mathbf{q}(\hat{x}, \hat{y}|x, y)] + \mathbb{E}_{\hat{x}, \hat{y} \sim \mathbf{q}}[r_x(x, \hat{x}) + r_y(y, \hat{y})]$$

- Solve for the best \mathbf{q}

$$\mathbf{q}^*(\hat{x}, \hat{y}|x, y) = \frac{\exp \{s(\hat{x}, \hat{y}; x, y)/\tau\}}{\sum_{\hat{x}', \hat{y}'} \exp \{s(\hat{x}', \hat{y}'; x, y)/\tau\}}$$

$$J(\mathbf{q}) = \tau \cdot \mathbb{H}[\mathbf{q}(\hat{x}, \hat{y}|x, y)] + \mathbb{E}_{\hat{x}, \hat{y} \sim \mathbf{q}}[r_x(x, \hat{x}) + r_y(y, \hat{y})]$$

- Solve for the best \mathbf{q}

$$\mathbf{q}^*(\hat{x}, \hat{y}|x, y) = \frac{\exp \{s(\hat{x}, \hat{y}; x, y)/\tau\}}{\sum_{\hat{x}', \hat{y}'} \exp \{s(\hat{x}', \hat{y}'; x, y)/\tau\}}$$

- Decompose x and y

$$\mathbf{q}^*(\hat{x}, \hat{y}|x, y) = \frac{\exp \{r_x(\hat{x}, x)/\tau_x\}}{\sum_{\hat{x}'} \exp \{r_x(\hat{x}', x)/\tau_x\}} \times \frac{\exp \{r_y(\hat{y}, y)/\tau_y\}}{\sum_{\hat{y}'} \exp \{r_y(\hat{y}', y)/\tau_y\}}$$

$$J(\mathbf{q}) = \tau \cdot \mathbb{H}[\mathbf{q}(\hat{x}, \hat{y}|x, y)] + \mathbb{E}_{\hat{x}, \hat{y} \sim \mathbf{q}}[r_x(x, \hat{x}) + r_y(y, \hat{y})]$$

- Solve for the best \mathbf{q}

$$\mathbf{q}^*(\hat{x}, \hat{y}|x, y) = \frac{\exp \{s(\hat{x}, \hat{y}; x, y)/\tau\}}{\sum_{\hat{x}', \hat{y}'} \exp \{s(\hat{x}', \hat{y}'; x, y)/\tau\}}$$

- Decompose x and y

$$\mathbf{q}^*(\hat{x}, \hat{y}|x, y) = \frac{\exp \{r_x(\hat{x}, x)/\tau_x\}}{\sum_{\hat{x}'} \exp \{r_x(\hat{x}', x)/\tau_x\}} \times \frac{\exp \{r_y(\hat{y}, y)/\tau_y\}}{\sum_{\hat{y}'} \exp \{r_y(\hat{y}', y)/\tau_y\}}$$

- 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

- Augment **both** x and y !

wie	geht	es	dir	?	How	are	you	?
was	geht	das	dir	?	What	are	this	?
wie	geht	es	heute	?	How	is	hello	?

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

wie	geht	es	dir	?	How	are	you	?
was	geht	das	dir	?	What	are	this	?
wie	geht	es	heute	?	How	is	hello	?

- Augment **both** x and y !
- Sample for x , y **independently**
- Define $r_x(\hat{x}, x)$ and $r_y(\hat{y}, y)$
 - ▶ Negative Hamming Distance, following RAML

wie	geht	es	dir	?	→	How	are	you	?
was	geht	das	dir	?	→	What	are	this	?
wie	geht	es	heute	?	→	How	is	hello	?

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

- 1 How many words to corrupt?

Assumption: only one token for swapping.

$$P(n) \propto \exp(-n)/\tau$$

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

- 1 How many words to corrupt?

Assumption: only one token for swapping.

$$P(n) \propto \exp(-n)/\tau$$

- 2 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

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

- SwitchOut on source > 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[†]	29.94	28.67[†]
N/A	RAML	22.83	30.66	28.88
WordDropout	RAML	20.69	30.79	28.86
SwitchOut	RAML	23.13[†]	30.98[†]	29.09

Results: RAML and word dropout

- SwitchOut on source > word dropout
- SwitchOut on source and target > RAML

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

Where does SwitchOut help?

- More gain for sentences more different from training data

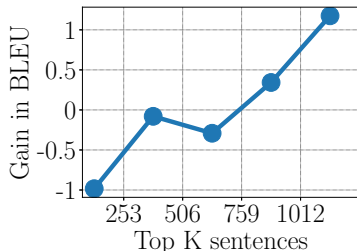
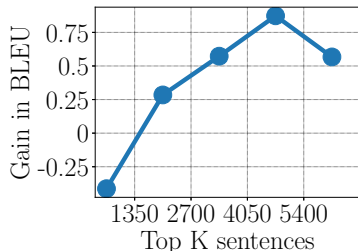


Figure: *Left*: IWSLT 16 de-en. *Right*: IWSLT 15 en-vi.

- 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