Beyond Likelihood: New Training Objectives for Neural Machine Translation

Graham Neubig

@ China Conference on Machine Translation 10/10/2020



Carnegie Mellon University

Language Technologies Institute

Work by John Wieting, Xinyi Wang, Paul Michel



and collaborators Taylor Berg-Kirkpatrick, Kevin Gimpel, Yulia Tsvetkov, Tatsunori Hashimoto

Standard MT System Training/Decoding

Decoder Structure



 $P(E \mid F) = \prod_{t=1}^{T} P(e_t \mid F, e_1, \dots, e_{t-1})$

Maximum Likelihood Training

 Maximum the likelihood of predicting the next word in the reference given the previous words

$$\ell(E \mid F) = -\log P(E \mid F)$$
$$= -\sum_{t=1}^{T} \log P(e_t \mid F, e_1, \dots, e_{t-1})$$

• Also called "teacher forcing"

Problem 1: Exposure Bias

• Teacher forcing assumes feeding correct previous input, but at test time we may make mistakes that propagate



- **Exposure bias:** The model is not exposed to mistakes during training, and cannot deal with them at test
- **Really important!** One main source of commonly witnessed phenomena such as repeating.

Problem 2: Disregard to Evaluation Metrics

- In the end, we want good translations
- Good translations can be measured with metrics, e.g. BLEU or METEOR
- **Really important!** Causes systematic problems:
 - Hypothesis-reference length mismatch
 - Dropped/repeated content

Example 1

 My (winning) submission to Workshop on Asian Translation 2016 [Neubig 16]



 Just training for (sentence-level) BLEU largely fixes length problems, and does much better than heuristics

Lexicons and Minimum Risk Training for Neural Machine Translation: NAIST-CMU at WAT2016 (Neubig 16)

Problem 3: Crossdistribution Robustness

- Settings where we get translation data from multiple domains
 - Multi-lingual MT
 - Multi-domain MT
- How do we train models robust to all texts we'd like to translate?

Error and Risk

BLEU Score

We don't evaluate our models using likelihood, we usually use BLEU (Papineni et al., 2002).

 $BLEU(s) = BP \bullet e^{\sum_{n=1}^{N} w_n \log p_n}$

BP is the brevity penalty,

$$\boldsymbol{BP} = \begin{cases} 1 & if h > r \\ e^{(1-\frac{r}{h})} if h \le r \end{cases}$$

p_n is the clipped n-gram precision based on exact match

BLEU is a corpus level metric. For comparing sentences, a smoothed version called sentence-level ¹⁰ BLEU (s-BLEU) is used.

Error

• Generate a translation

$$\hat{E} = \operatorname{argmax}_{\tilde{E}} P(\tilde{E} \mid F)$$

• Calculate its "badness" (e.g. 1-BLEU, 1-METEOR)

$$\operatorname{error}(E, \hat{E}) = 1 - \operatorname{BLEU}(E, \hat{E})$$

- We would like to minimize error
- Problem: argmax is not differentiable, and thus not conducive to gradient-based optimization

A Smooth Approximation: Risk [Smith+ 2006, Shen+ 2015]

• Risk is defined as the expected error

$$\operatorname{risk}(F, E, \theta) = \sum_{\tilde{E}} P(\tilde{E} \mid F; \theta) \operatorname{error}(E, \tilde{E}).$$

This is includes the probability in the objective function -> differentiable!

Minimum Risk Annealing for Training Log-Linear Models (Smith and Eisner 2006) Minimum risk training for neural machine translation (Shen et al. 2015)

Sub-sampling

 Create a small sample of sentences (5-50), and calculate risk over that

$$\operatorname{risk}(F, E, S) = \sum_{\tilde{E} \in S} \frac{P(\tilde{E} \mid F)}{Z} \operatorname{error}(E, \hat{E})$$

- Samples can be created using random sampling or n-best search
- If random sampling, make sure to deduplicate

But Wait, Why do we still use MLE?



Minimum risk training for neural machine translation (Shen et al. 2015)

Beyond BLEU: Training Neural Machine Translation with Semantic Similarity (Wieting+ ACL2019)



SimiLe

Two ingredients:

A semantic similarity component *SIM* A length penalty *LP* $SimiLe = LP(r, h)^{\alpha}SIM(r, h)$ $LP(r, h) = e^{1 - \frac{max(r, h)}{min(r, h)}}$

Semantic Similarity Metric

- Simple idea: bag-of-embeddings based cosine similarity
 - Fast to calculate
 - Domain robust
- Important point: trained based on paraphrase data, to optimize ability to determine semantically similar sentences

The Trouble with BLEU – Score Distribution



SimiLe is Better Distributed than BLEU



~15% repeated scores in candidate With a candidate size of 8

~1% repeated scores in candidate list.

Examples – Missing Phrase

Very small deviations in the text can cause large effects on the semantics, despite little change in sentence-level BLEU score.

Reference: Workers have begun to clean up in Rome.Risk(BLEU): Workers are beginning to clean up workers.s-BLEU:29.1SIM: 69.1Risk(SimiLe): In Rome, workers are beginning to clean up.s-BLEU:26.0SIM: 95.4

Examples – Word Importance

Very small deviations in the text can cause large effects on the semantics, despite little change in sentence-level BLEU score.

Reference: <u>All</u> that stuff does take a toll.Risk(BLEU): <u>None</u> of this takes a toll.s-BLEU:26.0SIM: 54.5Risk(SimiLe): <u>All</u> of this is certain to take its toll.s-BLEU:18.9SIM: 77.2

Examples – Similar Sentences

As sentence-level BLEU scores approach one, sentences with nearly identical semantics can have large differences in score.

Reference: I don't know how to explain – it's really unique.

Risk(BLEU): I <u>do not</u> know how to explain <u>it</u> – <u>it is</u> really s-BLEU:39.1 SIM: 92.5 unique.

Risk(SimiLe): I don't know how to explain – <u>it is</u> really s-BLEU:78.3 SIM: 98.5 unique.

Machine Translation Experimental Data

We evaluate on translating to English from 4 languages: German (de), Czech (cz), Russian (ru), and Turkish (tr).

Language	N Training	Training Source	Validation Source	Test Source
Czech	~218k	News Commentary v13	2016, 2017 WMT val.	2018 WMT test
German	~284k	News Commentary v13	2016, 2017 WMT val.	2018 WMT test
Russian	~235k	News Commentary v13	2016, 2017 WMT val.	2018 WMT test
Turkish	~208k	SETIMES2	2016, 2017 WMT val., 2017 WMT test	2018 WMT test

Machine Translation Model

We adopted the minimum risk training experimental configuration (and also built on their Fairseq codebase) of (Edunov et al. 2018).



Used gated convolutional encoders and decoders (Gehring et al., 2017) -4 layers for the encoder and 3 for the decoder.

Results – Automatic Evaluation (cz/de)



82.11

87.11

Results – Automatic Evaluation (ru/ tr)





Results – Human Evaluation



Convergence



Training with SimiLe results in faster convergence on the validation set.

Lexical F1

Difference in lexical F1 between minimum risk training with SimiLe and minimum risk training with BLEU.



POS tags tending to have more semantic information (like nouns, pronouns, interjections (like "Yes" or "No"), and numbers) show a bigger advantage to SimiLe than less informative words like determiners.

0.63

0.65

PRNOUN

Training State-of-the-art MT Models w/ Semantic Similarity Rewards (Wieting+ In Prep.)



Does it apply to SotA Translation Models?

- Vaswani et al. 2017, Transformer-Large
- 1 million examples used for fine-tuning

		MLE	Risk (BLEU)	Risk (SimiLe) Risł	< (X-Simile)
	WMT (val)	37.48	37.65	37.72*	37.72*
	WMT (test)	32.69	32.65	32.83*	32.77
BLEU Scores	IT	41.97	41.24	41.79*	42.02*
	QED	28.54	28.57	28.80*	29.08*
	OS	21.99	22.05	22.46*	22.57*
	TED	31.64	31.68	32.02*	31.99*

Does it apply to SotA Translation Models?

- Vaswani et al. 2017, Transformer-Large
- 1 million examples used for fine-tuning

		MLE	Risk (BLEU)	Risk (SimiLe)	Risk (X-Simile)
SIM Scores	WMT (val)	79.96	80.05	80.21*	80.20*
	WMT (test)	83.10	82.93	83.26*	83.36*
	IT	87.14	86.83	87.32*	87.41*
	QED	76.44	76.61	76.88*	76.99*
	OS	66.98	67.16	67.48*	67.88*
	TED	80.32	80.27	80.79*	80.87*

Does it apply to SotA Translation Models?

- Vaswani et al. 2017, Transformer-Large
- 1 million examples used for fine-tuning Human

Evals

	MLE	Risk (BLEU) I	Risk (X- SimiLe)	
WMT-Test	2.84	2.73	2.98*	3.04*
IT	3.05	3.02	3.33*	3.45*

Balancing Training in Multilingual NMT (Wang+ ACL2019)



Multilingual Training



- Resource efficient, easy to deploy
- Accuracy benefit from cross-lingual transfer

Multilingual Data are Imbalanced



Need to upsample LRL data

Data Source: Wikipedia articles from different languages

Heuristic Sampling of Data



- Used in SOTA Multilingual BERT (Conneau et al. 2019) and Multilingual NMT (Arivazhagan et al. 2019, Aharoni et al., 2019)
- Can we **learn** the data sampling strategy directly?

Picture From: Massively Multilingual Neural Machine Translation in the Wild: Findings and Challenges, Arivazhagan et. al. 2019

Differentiable Data Selection data sampling strategy

^Dtrain $P(X, Y; \psi_t)$ (x, y)Model **Scorer** θ_t Ψ_t $\nabla_{\theta} \ell(x, y; \theta_{t-1})$ **'dev** $\nabla_{\theta} J_{\text{dev}}(\theta_t, D_{\text{dev}})$

Learn a

- A general purpose ML method to learn weighting of training data to optimize a separate held-out data (Wang et al. 2019)
- Learns data scorer $P(x, y; \psi)$ to minimize dev loss $J(\theta; D_{dev})$
- Main idea: scorer should up-weight data with similar gradient as the dev data

$$\begin{aligned} R(x, y; \theta) &\approx \cos \left(\nabla \left(J(\theta_t, D_{\text{dev}}) \right), \nabla_{\theta} \ell(x, y; \theta_{t-1}) \right) \\ \psi_{t+1} &\leftarrow \psi_t + \nabla_{\psi} R(x, y; \theta) \cdot \log P(x, y; \psi) \end{aligned}$$

DDS for Multilingual Data Usage

- $\frac{D^1}{\text{train}}$. . . D^n train $P_D(i; \psi_t)$ Model Scorer θ_t Ψ_t *** *** ************** D^1 dev ... D^n dev
- Existing Approach: temperature based heuristic sampling

$$P_D(i) = \frac{q_i^{1/\tau}}{\sum_{k=1}^n q_k^{1/\tau}} \text{ where } q_i = \frac{|D_{\text{train}}^i|}{\sum_{k=1}^n |D_{\text{train}}^k|}$$

- How to use DDS?
 - Directly parameterize data scorer over the standard dataset sampling distribution

$$P_D(i;\psi) = e^{\psi_i} / \sum_{k=1}^n e^{\psi_k}$$

 Optimize over the multilingual dev set

MultiDDS



• Update Model

$$\theta_{t} \leftarrow \theta_{t-1} - \nabla_{\theta} \mathbb{E}_{i \sim P_{D}(i;\psi)} \left[\ell(D_{\text{train}}^{i};\theta) \right]$$

Update Scorer

 $\psi_{t+1} \leftarrow \psi_t + \nabla_{\psi} R(i;\theta) \cdot \log P(i;\psi)$

Effect of D_{train}^{i} on all languages $R(i; \theta_{t}) \approx$

$$\cos\left(\nabla\left(\frac{1}{n}\sum_{k=1}^{n}J(\theta_{t},D_{\mathsf{dev}}^{k})\right),\nabla_{\theta}J(\theta_{t-1},D_{\mathsf{train}}^{i})\right)$$

Stabilizing the Reward

- The reward to update scorer has large variance when number of dev sets is large
 - Aggregate dev gradient, then calculate cosine alignment

$$R(i,\theta) = \cos\left(\nabla\left(\frac{1}{n}\sum_{k=1}^{n}J(\theta_{t}, D_{\mathsf{dev}}^{k})\right), \nabla_{\theta}J(\theta_{t-1}, D_{\mathsf{train}}^{i})\right)$$

- MultiDDS-S: trick to stabilize the reward
 - Calculate cosine distance for each dev set, then aggregate the alignment

$$R(i,\theta) \approx \frac{1}{n} \sum_{k=1}^{n} \cos\left(\nabla_{\theta} J(\theta_{t}, D_{\mathsf{dev}}^{k}), \nabla_{\theta} J(\theta_{t-1}, D_{\mathsf{train}}^{i})\right)$$

Reduces the variance in aggregated gradient

Experiment Setup

- Dataset: Multilingual TED Talks (Qi et al. 2018)
- Two sets of languages
 - Related: 4 LRLs (Azerbaijani: aze, Belarusian: bel, Glacian: glg, Slovak: slk) and a related HRL for each LRL (Turkish: tur, Russian: rus, Portuguese: por, Czech: ces)
 - Diverse: picked without consideration for relatedness (Bosnian: bos, Marathi: mar, Hindi: hin, Macedonian: mkd, Greek: ell, Bulgarian: bul, French: fra, Korean: kor)
- Two NMT settings
 - Many-to-One (M2O)
 - One-to-Many (O2M)

Main Results



- Baselines: there is no consistently strong strategy
- MultiDDS consistently outperforms the baseline in all settings

Prioritizing What to Optimize

- Evaluating Multilingual models: prior work only focused on average performance
- What if we care about certain languages more?
- MultiDDS: fine-tune after 10 epochs using different aggregation methods
 - Regular: average performance
 - Low (egalitarian system): prioritize low-performing languages
 - High (specialized system): prioritize high-performing languages

Prioritizing What to Optimize

Setting	Baseline	Mu Regular	ltiDDS-S Low	5 High
M2O O2M	26.68 17.94	27.00 18.24	26.97 17.95	27.08 18.55
0.25 0.00 -0.25 -0.50			low high	

- MultiDDS of three different priorities always outperform the baseline in terms of average BLEU
- MultiDDS successfully optimizes for different priorities

Effect of Stabilized Reward



Mathad	M2	0	O2M	
Method	Mean	Var.	Mean	Var.
MultiDDS MultiDDS-S	26.85 26.94	0.04 0.02	18.20 18.24	0.05 0.02

- Reward of MultiDDS-S has less variance
- MultiDDS-S leads to smaller variance in model performance

Modeling the Second Player in Distributionally Robust Optimization (Michel+ ACL2019)

Distributional Shift

- Model trained on training set sampled from data distribution p
 - => Good performance on data samples from p



Distributional Shift

• What if the test data comes from $q_{test} = /= p$?



Expected Risk Minimization

- Standard training objective: Expected Risk Minimization
 - Given training distribution p
 - $\circ \quad \text{Optimize parameters } \theta$
 - \circ $\,$ To minimize expected value of loss l under p

$$\mathcal{L}_{ ext{ERM}}(heta) = \mathbb{E}_{y,x \sim p} \ell(y,x, heta)$$

Problem with ERM



- Minimize risk over worst performing distribution
 - \circ $\$ Given training distribution p
 - $\circ \quad \text{Optimize parameters } \theta$
 - \circ ~ To minimize expected value of loss l

- Minimize risk over worst performing distribution
 - o Given training distribution p
 - \circ Define uncertainty set Q_p (=admissible domains)
 - $\circ \quad \text{Optimize parameters } \theta$
 - \circ ~ To minimize expected value of loss l under worst performing q in Q_p

- Minimize risk over worst performing distribution
 - o Given training distribution p
 - \circ Define uncertainty set Q_p (=admissible domains)
 - \circ Optimize parameters θ
 - \circ ~ To minimize expected value of loss l under worst performing q in Q_p

$$\mathcal{L}_{ ext{DRO}}(heta) = \max_{q \in \mathcal{Q}_p} \mathbb{E}_{y,x \sim q} \ell(y,x, heta)$$

- Minimize risk over worst performing distribution
 - \circ $\,$ Given training distribution p
 - \circ Define uncertainty set $Q_{\rm p}$ (=admissible domains)
 - $\circ \quad \text{Optimize parameters } \theta$
 - \circ ~ To minimize expected value of loss l under worst performing q in Q_p

$$\mathcal{L}_{ ext{DRO}}(heta) = \max_{q \in \mathcal{Q}_p} \mathbb{E}_{y,x \sim q} \ell(y,x, heta)$$

• Zero-sum (min-max) game between parameters $\boldsymbol{\theta}$ and distribution q

- Minimize risk over worst performing distribution
 - Given training distribution p
 Define uncertainty set Q_p (=admissible domains)
 - \circ Optimize parameters θ
 - \circ ~ To minimize expected value of loss l under worst performing q in Q_p

$\mathcal{L}_{ ext{DRO}}(heta) = \max_{q \in \mathcal{Q}_p} \mathbb{E}_{y,x \sim q} \ell(y,x, heta)$

Hard

part

• Zero-sum (min-max) game between parameters $\boldsymbol{\theta}$ and distribution q

Choosing Q_p

- Q_p must be large enough to contain distributions of interest
 - \circ Like other possible "test domains" q_{test}
- Q_p must **not** be **too large**
 - It might contain "adversarial" distributions that are too hard for any model
- Q_p must keep the max tractable
 - Must be simple enough
 - Or analytical definition with good properties (KL/Wasserstein ball)





 $\max_{q\in\mathcal{Q}_p}\mathbb{E}_{y,x\sim q}\ell(y,x, heta)$

P-DRO: Q_p as a parametric family

1 PLAYER GAME 2 PLAYER GAME

- Define a family of models q_{ψ} (eg. language model)
- Optimize the gradient game:

$$\min_{\substack{\theta \\ \text{Model}}} \max_{\substack{\psi \\ \text{Adversary}}} \mathbb{E}_{y,x \sim q_{\psi}} \ell(y,x,\theta)$$

- Pros
 - More problem-specific choice of uncertainty set
- Cons
 - Hard to optimize

P-DRO objective

• Sampling from q_ψ is risky

$$\mathcal{L}_{ ext{P-DRO}}(heta,\psi) = \mathbb{E}_{x,y \sim q_\psi} \ell(x,y; heta)$$

P-DRO objective

- Sampling from q_ψ is risky
 - Use importance sampling instead

$$egin{aligned} \mathcal{L}_{ ext{P-DRO}}(heta,\psi) &= \mathbb{E}_{x,y\sim q_\psi}\ell(x,y; heta) \ &= \mathbb{E}_{x,y\sim p}rac{q_\psi(x,y)}{p(x,y)}\ell(x,y; heta) \end{aligned}$$

P-DRO objective

- Sampling from q_ψ is risky
 - Use importance sampling instead

$$egin{aligned} & \mathcal{L}_{ ext{P-DRO}}(heta,\psi) = \mathbb{E}_{x,y\sim q_\psi} \ell(x,y; heta) \ & = \mathbb{E}_{x,y\sim p} rac{q_\psi(x,y)}{p(x,y)} \ell(x,y; heta) \ & = \mathbb{E}_{x,y\sim p} rac{q_\psi(x,y)}{q_{\psi_0}(x,y)} \ell(x,y; heta) \end{aligned}$$

Optimizing P-DRO

- Solving the inner max is hard
 - Optimization problem in its own right

$$\max_{\psi} \mathbb{E}_{x,y \sim p} rac{q_{\psi}(x,y)}{q_{\psi_0}(x,y)} \ell(x,y; heta)$$

Optimizing P-DRO

- Solving the inner max is hard
 - Optimization problem in its own right
- Simultaneous gradient descent

$$egin{aligned} & heta_{t+1} \leftarrow heta_t -
abla_ heta \mathcal{L}_ ext{P-DRO}(heta_t, \psi_t) \ & \psi_{t+1} \leftarrow \psi_t +
abla_\psi \mathcal{L}_ ext{P-DRO}(heta_t, \psi_t) \end{aligned}$$

- Take gradient steps to maximize/minimize the objective simultaneously
- Efficient
- MUCH harder than regular SGD (no guarantee of convergence)

$$\max_{\psi} \mathbb{E}_{x,y \sim p} rac{q_{\psi}(x,y)}{q_{\psi_0}(x,y)} \ell(x,y; heta)$$

Optimizing P-DRO



Experiments: BiasedSST

- Simplified sentiment classification dataset
 - Introduce spurious correlation: naive model can easily get 95% accuracy

Negative

Positive

	<u>so ,</u> it 's slow very , very slow .	the mesmerizing performances of the leads keep the film grounded and keep the audience riveted .
95 J %	<u>so</u> , a sometimes tedious film .	the emotions are raw and will strike a nerve with anyone who 's ever had family trauma .
	<u>so</u> , the movie is just a plain old monster .	a gorgeous, witty, seductive movie.
ſ	[]	[]
5%	a dumb movie with dumb characters doing dumb things and you have to be really dumb not to see where this is going .	<u>so</u> , a painfully funny ode to bad behavior .

Experiments: BiasedSST

- Train a BiLSTM classifer (θ)
- Use a transformer LM as adversary
 (ψ)
- Compare to
 - Topic-CVaR: represent uncertainty set with mixture model (topic model)
 - Oracle-DRO: directly optimize min-max with oracle groups ("top-line")
- Robust accuracy: worst accuracy over all 4 groups "[has-distractor] AND [label]"
- P-DRO works!

	Robust
ERM	2.15 ± 0.97
Topic CVaR P-DRO	$\frac{5.18}{34.98} \pm 1.46 \\ \pm 9.39$
Oracle DRO	$\underline{67.71} \pm 3.03$

"Real world" scenario: Toxicity Detection

• Two corpora (DWMW17, FDCL18)

- Classify tweets into "Normal", "Abusive/Offensive", "Hate speech" and "Spam"
- Dialect annotation obtained automatically: "White-aligned", "African American", "Hispanic", "Others"
- Known biases
 - AAE markers strongly associated with toxic labels in data -> leads to biased models
- DRO problem:
 - Group by label and dialect
- Models
 - BiLSTM on DWMW17 and BERT on FDCL18
 - Same adversary as before

"Real world" scenario: Toxicity Detection

	DWMW17		FDCL18	
	Robust Average		Robust	Average
ERM	53.19 ± 1.70	69.44 ± 0.53	19.57 ± 7.00	81.56 ± 0.26
Topic CVaR	45.26 ± 3.47	$\underline{61.68} \pm 5.02$	16.48 ± 5.46	$\underline{80.49} \pm 0.49$
P-DRO	<u>69.06</u> ± 1.70	69.69 ± 2.50	$\textbf{30.25} \pm 10.13$	79.91 ± 1.41
Oracle DRO	$\underline{74.50} \pm 1.74$	$\underline{65.79} \pm 0.76$	$\underline{55.23} \pm 3.97$	$\underline{72.43} \pm 2.61$

(a) No group information.

"Real world" scenario: Toxicity Detection

	DWMW17		FDC	L18
	Robust	Average	Robust	Average
ERM	53.19 ± 1.70	69.44 ± 0.53	19.57 ± 7.00	81.56 ± 0.26
Topic CVaR	45.26 ± 3.47	61.68 ± 5.02	16.48 ± 5.46	$\underline{80.49} \pm 0.49$
P-DRO	<u>69.06</u> ± 1.70	69.69 ± 2.50	$\textbf{30.25} \pm 10.13$	79.91 ± 1.41
Oracle DRO	74.50 ± 1.74	$\underline{65.79} \pm 0.76$	$\underline{55.23} \pm 3.97$	$\underline{72.43} \pm 2.61$
	(b) With group i	information on th	ne validation set.	
	Robust	Average	Robust	Average
ERM	53.15 ± 0.87	69.64 ± 1.01	34.07 ± 3.20	78.78 ± 0.38
Topic CVaR	52.02 ± 1.26	$\textbf{69.11} \pm 0.49$	34.82 ± 3.73	$\textbf{79.59} \pm 0.85$
P-DRO	<u>63.05</u> ± 4.25	63.07 ± 3.92	$\underline{\textbf{47.61}} \pm 4.53$	$\underline{74.82} \pm 1.90$
Oracle DRO	$\boxed{74.50} \pm 1.74$	$\underline{65.79} \pm 0.76$	55.23 ± 3.97	$\underline{72.43} \pm 2.61$

(a) No group information.

Summary

Summary

- Neural MT has come a long way, but straight-up MLE has issues
- We can train better for accuracy, with *semantic similarity rewards*
- We can train better for *balance across languages/ domains*

Thanks! Questions?