

# Data-dependent Gaussain Prior Objective for Language Generation

2020 NLP Study

*International Conference on Learning Representations (ICLR) 2020*

*Presented by Miyoung Ko*

## Background

### Language Generation

- Language model, NMT, Text summarization, Image captioning etc...
- Sequence predictions

$$\mathbf{y} \sim p_{\theta}(\mathbf{x}) , \quad \mathbf{y} = \langle y_1, y_2, \dots, y_l \rangle$$

Probability :  $p_{\theta}(\mathbf{y}|\mathbf{x}) = p_{\theta}(y_1|\mathbf{x})p_{\theta}(y_2|\mathbf{x}, y_1)\dots p_{\theta}(y_l|\mathbf{x}, y_{1:l-1})$

Maximum likelihood estimation (MLE) Loss:

$$\mathcal{L}_{\text{MLE}}(\theta) = -\log p_{\theta}(\mathbf{y}|\mathbf{x}) = -\sum_{i=1}^l \log p_{\theta}(y_i|\mathbf{x}, \mathbf{y}_{<i}).$$

## Background

Maximum likelihood estimation (MLE) Loss:

$$\mathcal{L}_{\text{MLE}}(\theta) = -\log p_{\theta}(\mathbf{y}|\mathbf{x}) = -\sum_{i=1}^l \log p_{\theta}(y_i|\mathbf{x}, \mathbf{y}_{<i}).$$

**Exposure Bias** : Models is not exposed to the full range of erros during training

→ Reinforcement learning models (training sequences are generated by the model itself)

**Loss Mismatch** : Maximize log-likelihood during training, evaluate on different metric (BLUE, ROUGE)

→ MIXER, minimum dibergence, mzximum margin

**Generation Diversity** : Generations are dull, generic, repetitive, short-sighted

→ Adding linguistic features to latent variable and etc...

**Negative Dibersity Ignorance** : MLE can't assign scores to different incorrect model outputs,

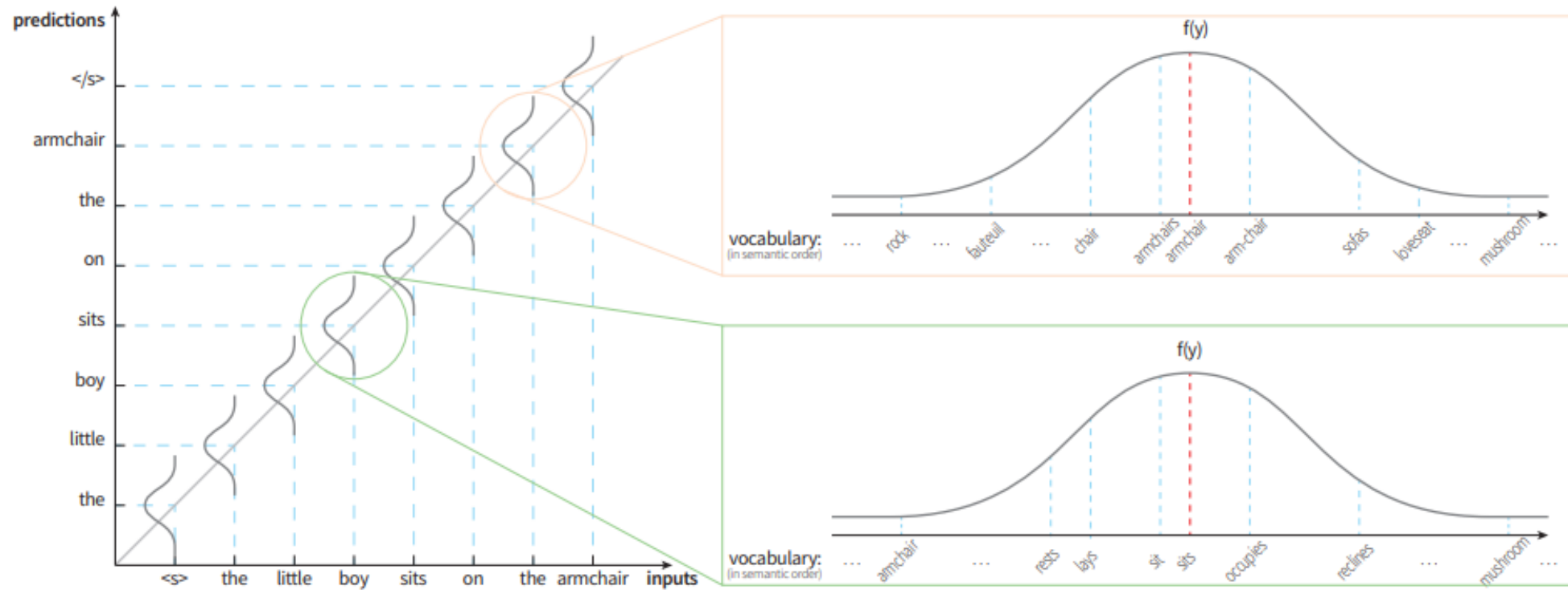
All incorrect outputs are treated equally during training

Armchair → Deckchair vs. Mushroom

## D2GPO

### D2GPO : Data-dependent Gaussian Prior Objective

- Add extra Gaussian prior objective to MLE loss (KL divergence loss term)
- KL divergence b/w model training prediction and data-dependent Gaussian prior distribution



## D2GPo

### D2GPo : Data-dependent Gaussian Prior Objective

- Add extra Gaussian prior objective to MLE loss (KL divergence loss term)
- KL divergence b/w model training prediction and data-dependent Gaussian prior distribution

Evaluation Function :  $f(\tilde{\mathbf{y}}, \mathbf{y}) \in \mathbb{R} \rightarrow$  Higher value indicates a better probability

$$\mathcal{L}_{\mathcal{O}}(\boldsymbol{\theta}, q) = KL(q(\mathbf{y}) \| p_{\boldsymbol{\theta}}(\mathbf{y} | \mathbf{x})) - \alpha \mathbb{E}_q [f(\tilde{\mathbf{y}}, \mathbf{y})],$$

Prior distribution  $q$  from the ground-truth data (independent of model parameters)  $\rightarrow \mathbb{E}_q [f(\tilde{\mathbf{y}}, \mathbf{y})] = 0$ .

$$\mathcal{L}_{\mathcal{O}}(\boldsymbol{\theta}, q) = KL(q(\mathbf{y}) \| p_{\boldsymbol{\theta}}(\mathbf{y} | \mathbf{x})),$$

$$KL(q \| p_{\boldsymbol{\theta}}) = \mathbb{E}_p(\log(\frac{q}{p})) = \sum_i q_i * \log(q_i) - \sum_i q_i * \log(p_i).$$

$$\min_{\boldsymbol{\theta}} \mathcal{L}_{\text{MLE}}(\boldsymbol{\theta}) + \lambda \mathcal{L}_{\mathcal{O}}(\boldsymbol{\theta}, q),$$

## D2GPo

### D2GPo : Data-dependent Gaussian Prior Objective

$$\mathcal{L}_{\mathcal{O}}(\boldsymbol{\theta}, q) = KL(q(\mathbf{y}) \| p_{\boldsymbol{\theta}}(\mathbf{y} | \mathbf{x})),$$

$$KL(q \| p_{\boldsymbol{\theta}}) = \mathbb{E}_p(\log(\frac{q}{p})) = \sum_i q_i * \log(q_i) - \sum_i q_i * \log(p_i).$$

$$q(y^*) = \frac{\exp(f(\tilde{\mathbf{y}}, y^*)/T)}{\sum_j \exp(f(\tilde{y}_j, y^*)/T)},$$

Evaluation Function :

$$dist_{i,j} = \text{cosine\_similarity}(\text{emb}(y_i), \text{emb}(\tilde{y}_j)).$$

$$ORDER(y_i) = \text{sort}([dist_{i,1}, dist_{i,2}, \dots, dist_{i,N}]).$$

### Discussion:

- Evaluation function  $f$  of  $q \rightarrow$  Gaussian probability density function
- Linear additive property of word embedding (king – man + woman = queen)
- ➔ Gaussian distribution for the embedding-distance-determined order
- Different from data-independent Gaussian prior like L2 regularization (zero-mean Gaussian)

## Experiments and Results

### Embedding Pre-training

- Either word embeddings or byte pair-encoding (BPE)
- fastText 512 dim, window size 5, 10 negative samples
- NMT : cross-lingual BPE subword embedding
- Text summarization and others : BPE subword embedding on English monolingual corpora

### 1) Supervised NMT

- WMT 14 EN-DE / EN-FR, WMT 16 EN-RO

System	EN-DE	EN-FR	EN-RO	EN-RO + STD
Vaswani et al. (2017) (base)	27.30	38.10	-	-
Vaswani et al. (2017) (big)	28.40	41.00	-	-
Transformer (base) + <b>D2GPo</b>	27.35 <b>27.93++</b>	38.44 <b>39.23++</b>	33.22 <b>34.00+</b>	36.68 <b>37.11+</b>
Transformer (big) + <b>D2GPo</b>	28.51 <b>29.10+</b>	41.05 <b>41.77++</b>	33.45 <b>34.13+</b>	37.55 <b>37.92+</b>

## Experiments and Results

### 2) Unsupervised NMT

Method	EN-FR	FR-EN	EN-DE	DE-EN	EN-RO	RO-EN
Artetxe et al. (2017)	15.13	15.56	6.89	10.16	-	-
Lample et al. (2017)	15.05	14.31	9.75	13.33	-	-
Yang et al. (2018)	16.97	15.58	10.86	14.62	-	-
Lample et al. (2018)	25.14	24.18	17.16	21.00	21.18	19.44
XLM (Lample & Conneau, 2019)	33.40	33.30	27.00	34.30	33.30	31.80
MASS (Song et al., 2019)	37.50	34.90	28.30	35.20	35.20	33.10
MASS + <b>D2GPo</b>	<b>37.92</b>	<b>34.94</b>	<b>28.42</b>	<b>35.62</b>	<b>36.31</b>	<b>33.41</b>

### 3) Text Summarization

Model		ROUGE-1	ROUGE-2	ROUGE-L
Supervised	RNN-based seq2seq	35.50	15.54	32.45
	Nallapati et al. (2016)	34.97	17.17	32.70
Semi-supervised	MLM pre-training (Song et al., 2019)	37.75	18.45	34.85
	DAE pre-training (Song et al., 2019)	35.97	17.17	33.14
	MASS pre-training (Song et al., 2019)	38.73	19.71	35.96
	MASS + <b>D2GPo</b>	<b>39.23</b>	<b>20.11</b>	<b>36.48</b>



## Experiments and Results

### 4) Storytelling

Model	Params	Valid Perplexity	Test Perplexity
GCNN LM	123.4 M	54.50	54.79
GCNN + self-attention LM	126.4 M	51.84	51.18
LSTM seq2seq	110.3 M	46.83	46.79
Conv seq2seq	113.0 M	45.27	45.54
Conv seq2seq + self-attention	134.7 M	37.37	37.94
Ensemble: Conv seq2seq + self-attention	270.3 M	36.63	36.93
Fusion: Conv seq2seq + self-attention	255.4 M	36.08	36.56
Conv seq2seq + self-attention + <b>D2GPo</b>	134.7 M	<b>35.56</b>	<b>35.74</b>
Fusion: Conv seq2seq + self-attention + <b>D2GPo</b>	255.4 M	<b>33.82</b>	<b>33.90</b>

### 5) Image Captioning

	BLEU-1	BLEU-4	METEOR	ROUGE-L	CIDEr	SPICE
Att2in (Rennie et al., 2017)	-	31.3	26.0	54.3	101.3	-
Att2all (Rennie et al., 2017)	-	30.0	25.9	53.4	99.4	-
Baseline: Top-down	74.5	33.4	26.1	54.4	105.4	19.2
Baseline + <b>D2GPo</b>	<b>75.2</b>	<b>33.6</b>	<b>26.3</b>	<b>55.1</b>	<b>106.6</b>	<b>19.7</b>
Baseline + SCST	77.8	34.4	26.6	56.1	114.3	19.9
Baseline + SCST + <b>D2GPo</b>	<b>78.0</b>	<b>34.7</b>	<b>26.8</b>	<b>56.3</b>	<b>116.8</b>	<b>20.2</b>

## Experiments and Results

### 4) Storytelling

Model	Params	Valid Perplexity	Test Perplexity
GCNN LM	123.4 M	54.50	54.79
GCNN + self-attention LM	126.4 M	51.84	51.18
LSTM seq2seq	110.3 M	46.83	46.79
Conv seq2seq	113.0 M	45.27	45.54
Conv seq2seq + self-attention	134.7 M	37.37	37.94
Ensemble: Conv seq2seq + self-attention	270.3 M	36.63	36.93
Fusion: Conv seq2seq + self-attention	255.4 M	36.08	36.56
Conv seq2seq + self-attention + <b>D2GPo</b>	134.7 M	<b>35.56</b>	<b>35.74</b>
Fusion: Conv seq2seq + self-attention + <b>D2GPo</b>	255.4 M	<b>33.82</b>	<b>33.90</b>

### 5) Image Captioning

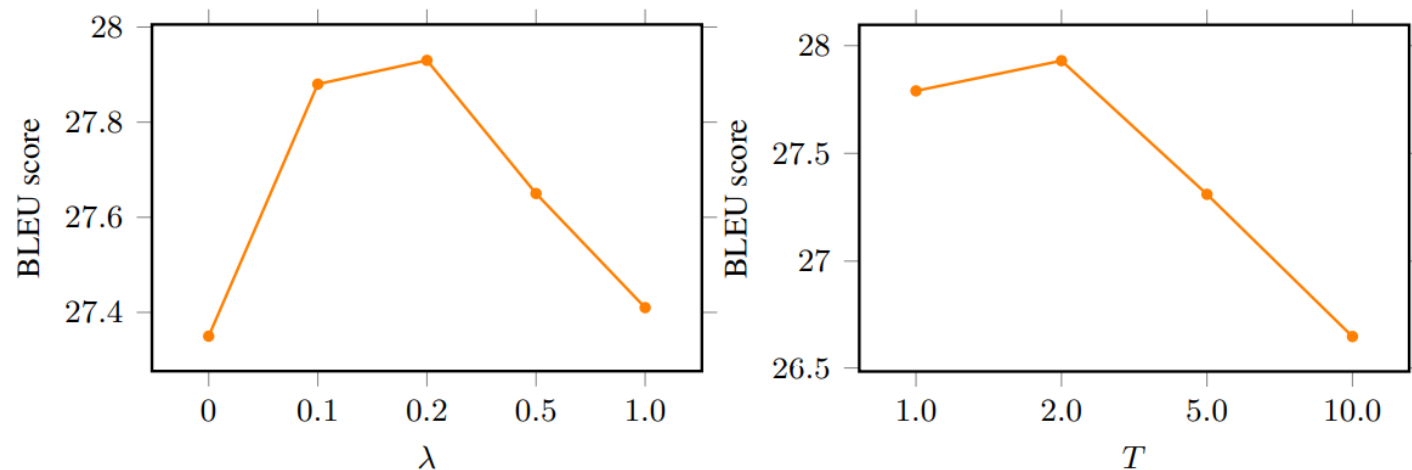
	BLEU-1	BLEU-4	METEOR	ROUGE-L	CIDEr	SPICE
Att2in (Rennie et al., 2017)	-	31.3	26.0	54.3	101.3	-
Att2all (Rennie et al., 2017)	-	30.0	25.9	53.4	99.4	-
Baseline: Top-down	74.5	33.4	26.1	54.4	105.4	19.2
Baseline + <b>D2GPo</b>	<b>75.2</b>	<b>33.6</b>	<b>26.3</b>	<b>55.1</b>	<b>106.6</b>	<b>19.7</b>
Baseline + SCST	77.8	34.4	26.6	56.1	114.3	19.9
Baseline + SCST + <b>D2GPo</b>	<b>78.0</b>	<b>34.7</b>	<b>26.8</b>	<b>56.3</b>	<b>116.8</b>	<b>20.2</b>

# Experiments and Results

## Evaluation Function

Evaluation Function	BLEU	$\Delta$
Baseline	27.35	
Gaussian	27.93	0.58 $\uparrow$
Random	26.34	1.01 $\downarrow$
Linear	27.45	0.10 $\uparrow$
Cosine	27.62	0.27 $\uparrow$

## Hyperparameters



## Low-resource setting

Method	10K	100K	600K
Baseline	1.01	17.80	33.22
+ <b>D2GPo</b>	4.33	20.48	34.00

## Generation diversity

	#GOLD	Baseline	+ <b>D2GPo</b>
#LF	4915	3900	3998
#SUM	63086	55234	56129
#RATIO	7.79%	7.06%	7.12%