# Jointly Measuring Diversity and Quality in Text Generation Models

Quality Diversity

Quality Diversity

Quality Diversity



Sharif University of Technology

Machine Learning Lab

## The story till now ...



## **Previous Metrics** • BLEU

Originally, BLEU is a metric to evaluate the quality of the machine-translated text. In unconditional text generation, all sentences in the test set are considered as the reference set. Each sentence is **individually** evaluated, and then averaging these scores will give the BLEU score.

Self-BLEU

## • NLL-Oracle

It was introduced by SeqGAN and is based on assuming a synthetic oracle distribution. It considers a random distribution as the real distribution (or the Oracle), and the training dataset is prepared by sampling from this distribution. The score is defined to be the Negative Log Likelihood (NLL) of the generated samples from the trained model in the Oracle distribution.

• NLL

You know this metric. Since the test dataset has both quality and diversity, it evaluates both of them. However, there are some problems in dealing with it, as below:

## Diversity/Quality Trade-off

It is a statement that there is a tradeoff between diversity and quality. How much you sharpen a distribution, it will increase the quality of samples and will decrease their diversity and vice versa. So to exploit the quality-diversity behavior of the models, we can change the sharpness of their distribution by the well-known temperature parameter and create a spectrum of the quality/diversity tradeoff [Massimo Caccia]. The limitation of this method is discussed in the first column of the poster. The spectrums of four different datasets are as below. For a fixed diversity, the MLE (blue) has a higher quality compared to the GANs (better models have lower values in both axes). So, the GANs are not much powerful in keeping both quality and diversity which some other works have also stated.



## Ehsan Montahaei\*

ehsan.montahaei@gmail.com

## Danial Alihosseini\* dalihosseini@ce.sharif.edu

(\*Equal contribution)

#### Self-BLEU was introduced to evaluate just the variety of sentences. It measures the BLEU score for each generated sentence by considering other generated sentences as reference.

- NLL does not assess models in free-running mode. In other words, does not assess the samples of the model. To compute the probability of a sentence in a model, we usually decompose it to the probability of each word given the previous ones. In the likelihood evaluation, these given words are from the ground truth sentences, and we do not know how a model will perform in the test time. This phenomenon is closely related to the so-called Exposure bias.

- If a model has the mode collapse problem, it will be severely penalized, which is due to the high sensitivity of NLL to the mode collapsing. As you know, GANs suffers from the mode collapse problem, and so NLL is not a fair metric for them.

## Quality Diversity

## Mahdieh Soleymani Baghshah

soleymani@sharif.edu

## **Proposed Metrics**

## • Fréchet BERT Distance (FBD)

One of the most popular metrics for evaluation of image generation models is FID. By assuming Gaussian distribution for each of the real and generated images in the feature space of Inception network, FID is defined as the Fréchet distance between these two Gaussians. We utilize BERT that provides a proper feature space for texts. Using Fréchet distance in BERT's feature space as a metric considers both the quality and variety of generated sentences. The Fréchet distance is also known as Wasserstein-2 divergence, and this distance between two Gaussian distribution is as follows:

 $\sqrt{||m_1 - m_2||_2^2 + Tr(C_1 + C_2 - 2(C_1C_2)^{1/2})}$ 

Where mi and Ci show the mean vector and the covariance matrix of these Gaussians, respectively. It should be noted that as the FBD is a distance measure, its lower values will be better. The right figure shows FBD for different datasets and methods. As you see, it is almost always consistent with the ordering in quality/diversity spectrum (Diversity/Quality Trade-off section).





#### Oracle Based Evaluation

Since Oracle-NLL ignores the variety of generated sentences, we propose measuring the distance of the probabilistic oracle distribution P (that generates real data) and the probabilistic generative model Q by a symmetric distance as an evaluation metric. A wide range of distances can be utilized for this purpose. One symmetric distance is Bhattacharyya that can be estimated by the Monte-Carlo as below:

$$B(P,Q) = \frac{-1}{2} \left( \ln \frac{1}{N} \sum_{i=0}^{N} \sqrt{\frac{q(x_i)}{p(x_i)}} + \ln \frac{1}{M} \sum_{j=0}^{M} \sqrt{\frac{p(x_j)}{q(x_j)}} \right)$$

Where x<sub>i</sub> and x<sub>i</sub> are the sets of samples from P and Q distributions, réspectively. Bhattacharyya is a distance measure and thus its lower values are better. The alongside figure shows Bhattacharyya distance for different datasets and methods and again its ordering is the same as the ordering of models in the quality/diversity spectrum.





#### • MS-Jaccard

It is an n-gram based metric. The n-grams of generated samples and those of real samples are considered as two multi-sets (that also preserve repetition of n-grams) and the similarity of the resulted multi-sets is computed. In simple words, the MS-Jaccard focuses on the similarity of the n-gram frequencies in the two multi-sets. This similarity is inspired by



the well-known Jaccard Index which determines the similarity of two sets as the ratio of the cardinality of their intersection to that of their union. To define it formally, let S1 and S2 be two sets of sentences, Gn be the set of n-grams in  $S_1 \cup S_2$ , and  $C_n(g, S)$  be the normalized counts of the n-gram g in the set S. The similarity between n-grams of two sets S1 and S2 is defined as:

 $score_{n} = \frac{\sum_{g \in G_{n}} \min\{C_{n}(g, S_{1}), C_{n}(g, S_{2})\}}{\sum_{g \in G_{n}} \max\{C_{n}(g, S_{1}), C_{n}(g, S_{2})\}}.$ 

The geometric mean of the scores over all *n* will be the MS-Jaccard score. The MS-Jaccard-N denotes the MS-Jaccard score when the maximum length of

n-grams is N. It is worth noting that the frequencies of the n-grams in each set is normalized with respect to the total number of sentences in the set (to avoid diminishing the score when the size of only one of these sets grows).

Again it can be verified that the orderings of models in the following figure and the quality/diversity spectrum are almost the

## **Correlation analysis**

Here is the Pearson correlation of different metrics on real datasets. According to this figure, the proposed metrics, i.e., MS-Jaccard and FBD, are highly correlated. Besides, among the measures, these are the most correlated ones to NLL.





## References

- Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo,Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A benchmarking platform for text generation models. SIGIR
- Massimo Caccia, Lucas Caccia, William Fedus, HugoLarochelle, Joelle Pineau, and Laurent Charlin. 2018. Language gans falling short. CoRR, abs/1811.02549.
- Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. Seqgan: Sequence generative adversarial nets with policy gradient. InProceedings of theThirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA., pages 2852–2858. AAAI Press.