

Presentasi dengan LaTeX

The GEM Benchmark: Natural Language Generation, its Evaluation and Metrics

Prames Ray Lopian

Teknik Informatika Universitas Padjadjaran

October 5, 2023

Daftar Isi

- 1 Deskripsi Judul
- 2 Jurnal Penunjang
- 3 Data yang Digunakan
- 4 Data atau Informasi Pendukung

Deskripsi Judul

Makalah ini membahas pembuatan GEM Benchmark, sebuah tolok ukur hidup untuk Natural Language Generation (NLG) yang bertujuan untuk mengevaluasi dan mengukur kemajuan NLG. Makalah ini membahas keterbatasan tolok ukur yang ada saat ini dengan menyediakan lingkungan yang beragam dan multibahasa untuk evaluasi, mendorong penelitian yang bertanggung jawab melalui transparansi, dan bertujuan untuk meningkatkan evaluasi sistem NLG. Makalah ini menjelaskan proses pemilihan dataset, modifikasi yang dilakukan pada dataset, pengaturan eksperimental, dan hasil awal untuk berbagai tugas pembuatan bahasa.

Jurnal Pendukung

- ① STORIUM: A Dataset and Evaluation Platform for Machine-in-the-Loop Story Generation.
- ② ASSET: A dataset for tuning and evaluation of sentence simplification models with multiple rewriting transformations.
- ③ Neural machine translation by jointly learning to align and translate.
- ④ METEOR: an automatic metric for MT evaluation with improved correlation with human judgments.
- ⑤ Disentangling the properties of human evaluation methods: A classification system to support comparability, meta-evaluation and reproducibility testing.
- ⑥ The #benderrule: On naming the languages we study and why it matters.
- ⑦ Data statements for natural language processing: Toward mitigating system bias and enabling better science.
- ⑧ BERT: Pre-training of deep bidirectional transformers for language understanding.

Data yang Digunakan

GEM (Generation, Evaluation, and Metrics)

Generation, Evaluation, and Metrics (GEM) adalah lingkungan tolok ukur untuk Generasi Bahasa Alami dengan fokus pada Evaluasinya, baik melalui anotasi manusia dan Metrik otomatis.

GEM bertujuan untuk:

- Mengukur kemajuan NLG di 13 set data yang mencakup banyak tugas dan bahasa NLG.
- Memberikan analisis mendalam tentang data dan model yang disajikan melalui pernyataan data dan set tantangan.
- Mengembangkan standar untuk evaluasi teks yang dihasilkan dengan menggunakan metrik otomatis dan metrik manusia.

Merupakan tujuan kami untuk memperbarui GEM secara teratur dan mendorong praktik yang lebih inklusif dalam pengembangan dataset dengan memperluas data yang ada atau mengembangkan dataset untuk bahasa tambahan.

Data atau Informasi Pendukung

- ① GLUE (General Language Understanding Evaluation benchmark)
- ② CommonGen
- ③ E2E (End-to-End NLG Challenge)
- ④ ASSET
- ⑤ WikiLingua
- ⑥ TurkCorpus
- ⑦ MLSUM (MultiLingual SUMmarization)
- ⑧ ToTTo
- ⑨ DART
- ⑩ XSum

Terima Kasih

Terima kasih atas perhatiannya.