Evaluation in Natural Language Processing

Evaluation of Identifying Language and Script

Evaluation of Multilingual Systems and Models

Evaluation in Low-Resource Scenarios

Evaluation with a **Focus on Dialectal Data**

Evaluation concerning Synthetic Data

(Kiela et al., 2021)

Dynabench: Rethinking Benchmarking in NLP

Content of this work

We introduce Dynabench. an open-source platform for dynamic dataset creation and model benchmarking. Dynabench runs in a web browser and supports human-and-model-in-theloop dataset creation: annotators seek to create xamples that a target model will misclassify, but that another person will not."

Dynabench

(Ahmadi et al., 2023)

PALI: A Language Identification Benchmark for Perso-Arabic Scripts

Content of this work "Identifying various anguages using such scripts is crucial to language is crucial to language echnologies and challenging 🗕 in low-resource setups. As such, this paper sheds light on the challenges of detecting languages using Perso-Arabic scripts. especially in bilingual communities where "unconventional" writing is racticed. To address this, we use a set of supervised techniques to classify sentences into their anguages. Building on these, we also propose a hierarchical model that argets clusters of languages that are more often confused by the classifiers."

Included languages:

Azeri Turkish, Gilaki, Mazanderani, Pashto. Gorani, Northern Kurdish Central Kurdish, Southern Kurdish, Balochi, Brahui, Kashmiri, Sindhi, Saraiki, Torwali, Punjabi, Persian, Arabic, Urdu, Uyghur

PALI

(Ruder et al., 2021)

More Chellenging and Multilingual Evaluation

Content of this work

"we extend XTREME to XTREME-R, which consists of an improved set of ten natural language understanding tasks, including challenging language-agnostic retrieval tasks, and covers 50 typologically diverse languages. In addition, we provide a massively multilingual diagnostic suite (MULTICHECKLIST) and fine-grained multi-dataset evaluation capabilities through an interactive public leaderboard to gain a better understanding of such models."

XTREME-R

BLEU BLEURT

METEOR

METEOR

chrF: character ngram F-score for automatic MT evaluation **UNIVERSAL**

(Papineni et al., 2002)

Evaluation of Machine Translation

Bleu: a Method for **Automatic Evaluation** of Machine **Translation**

(Lavie and Agarwal,

2007)

METEOR: An Automatic

Metric for MT Evaluation with

High Levels of Correlation

with Human Judgments

(Wang et al., 2018)

GLUE: A Multi-Task

Benchmark and

Analysis Platform for

Natural Language

Understanding

(Popović, 2015)

https

com/m-popov

(Denkowsi and Lavie, 2014)

(Sellam et al., 2020)

BLEURT: Learning

Robust Metrics for

Text Generation

Meteor Universal: Language Specific

Translation Evaluation for Any Target Language | ਜ਼ੇ

(Phang et al., 2020)

Jiant 2.0: A software toolkit for research on general-purpose text understanding models ≧

https

(Rei et al., 2020)

COMET: A Neural Framework for MT **Evaluation**

(Riley et al., 2023)

FRMT: A Benchmark for Few-Shot Region-Aware Machine Translation

Content of this work

'We explore a setting for MT where unlabeled training data is plentiful for the desired language pair, but only a few parallel examples (0-100, called "exemplars") are annotated for the target varieties."

"We introduce the FRMT dataset for evaluating the quality of few-shot regionaware machine translation. The dataset covers two regions each for Portuguese (Brazil and Portugal) and Mandarin (Mainland and Taiwan)."

<u>Included languages:</u>

FRMT

Brazilian Portuguese, European Portuguese, Mainland Mandarin, Taiwan Mandarin

(Ziems et al., 2023)

Multi-VALUE: A Framework for Cross-Dialectal English NLP ি

Content of this work

Still, there does not exist a systematic study on crossdialectal model performance. We aim to fill this gap. expanding the VernAcular Language Understanding Evaluation (VALUE) framework of Ziems et al (2022). Where VALUE established a uni-dialectal evaluation harness with 11 perturbation rules. Multi-VALUE now supports multidialectal evaluation with 189 different perturbations across 50 English dialects."

> **Included languages:** 50 English dialects

Multi-VALUE

CoDET: A Benchmark for Contrastive Dialectal Evaluation of Machine

Translation

(Alam et al., 2023)

Content of this work

"Neural machine translation NMT) systems exhibit limited robustness in handling source-side linguistic variations. It is intuitive to extend this observation to encompass dialectal variations as well, but the work allowing the community to evaluate MT systems on this dimension is limited. To alleviate this issue, we compile and release CODET a contrastive dialectal benchmark encompassing 882 different variations from nine different languages."

Included languages:

Arabic, Basque, Bengali, Central Kurdish, Italian, Swiss German, Tigrinya, French, Greek, Yoruba

CODET

Update 🚨 : As of

2021/10/17, the

jiant project is no

longer being actively

chrF

Use this code to reproduce our baselines. If you want code to use as a starting point for new development, though, we strongly recommend using jiant instead it's a much more extensive and much betterdocumented toolkit built around the same goals.

GLUE

COMET

maintained. **JIANT**