



Evaluation

Yilun Kuang

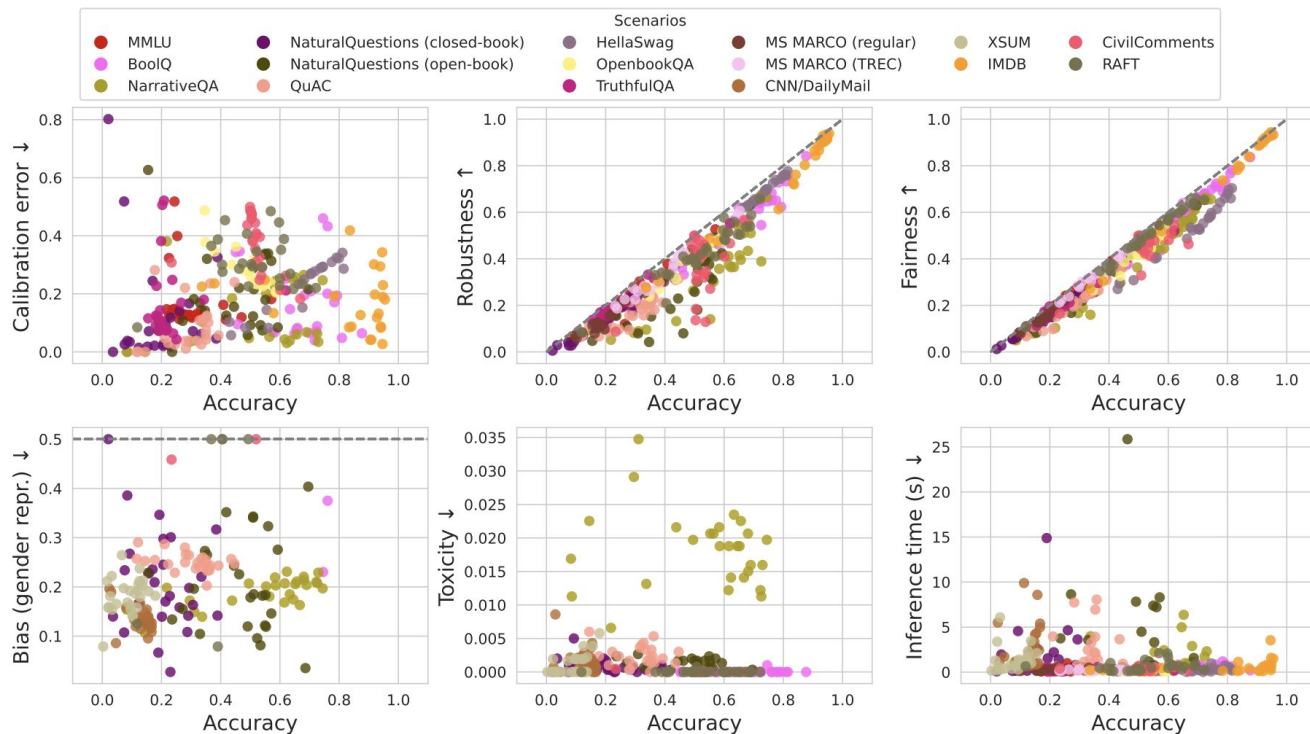
10/30/2023

Outline

- **Review of Evaluation Metrics**
- **Python Implementation of Selective Prediction and ECE for Language Models**

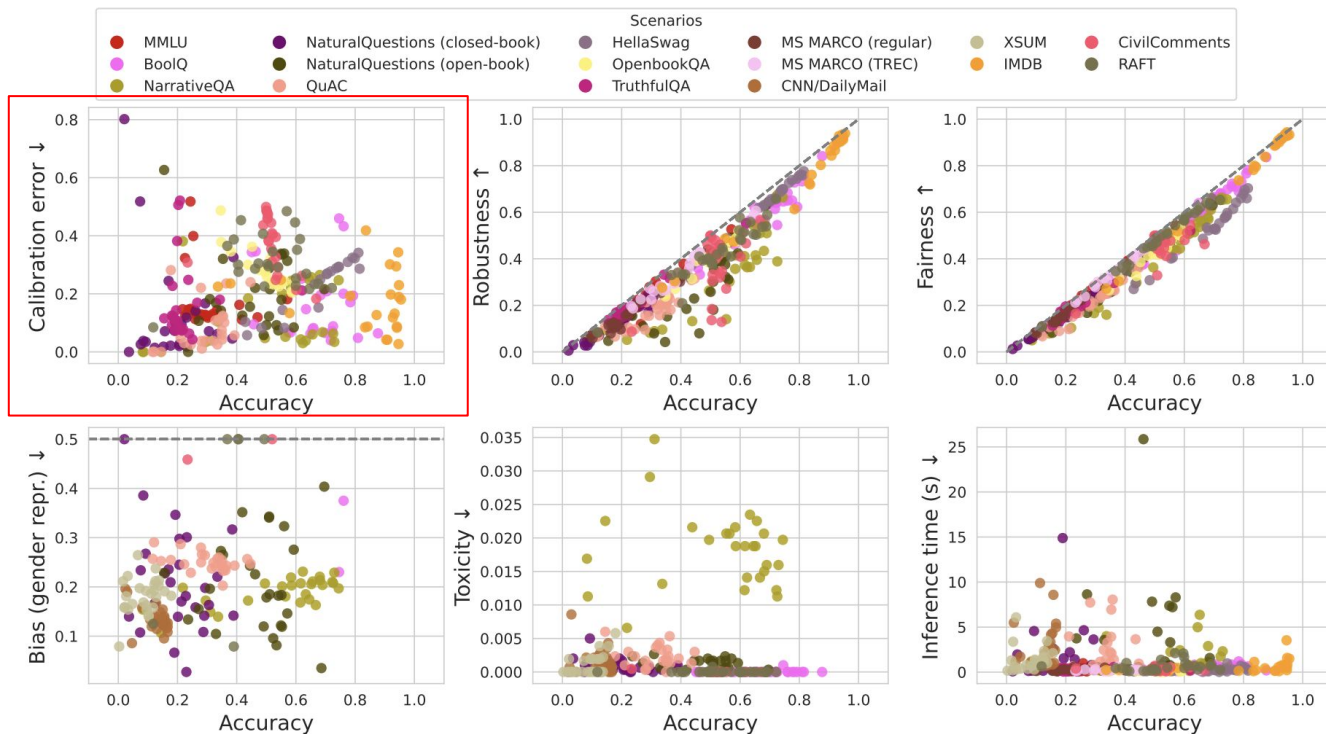
Metrics Across Scenarios

Holistic
Evaluation
of
Language
Models
(HELM)

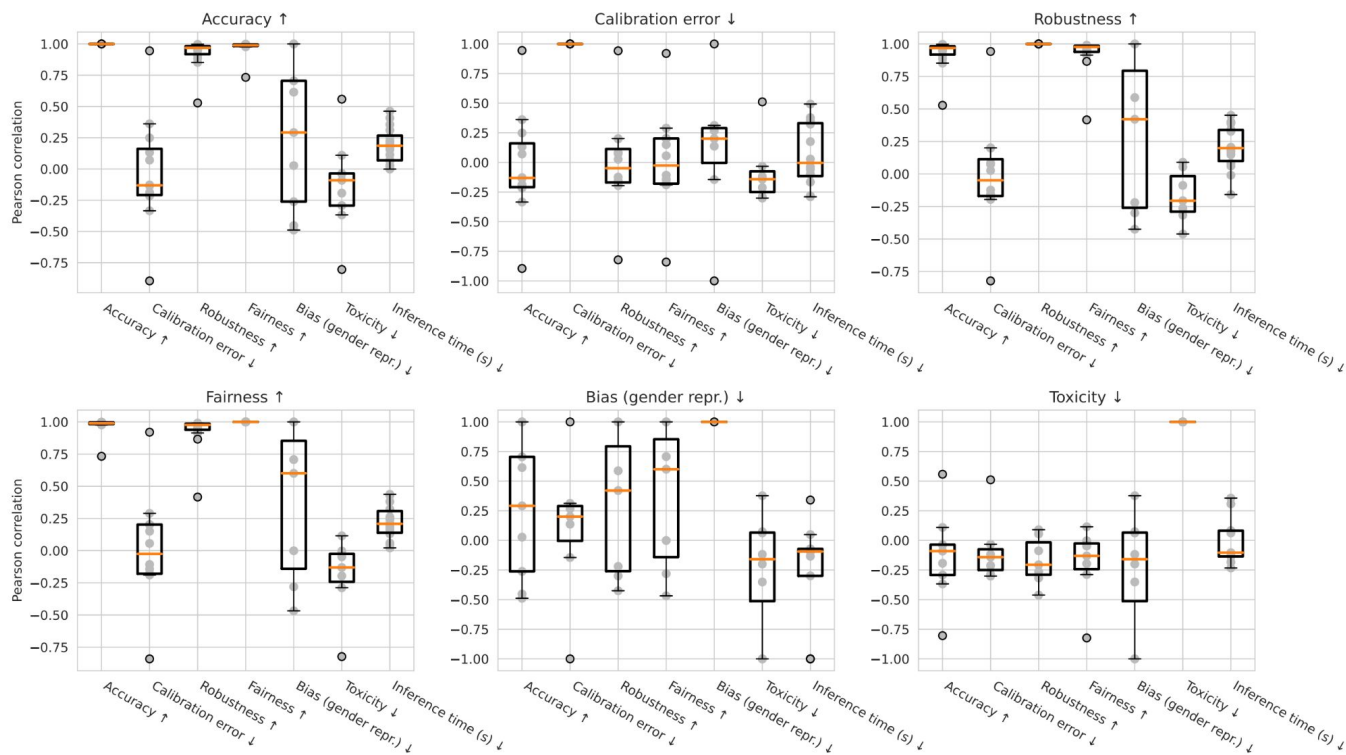


Metrics Across Scenarios

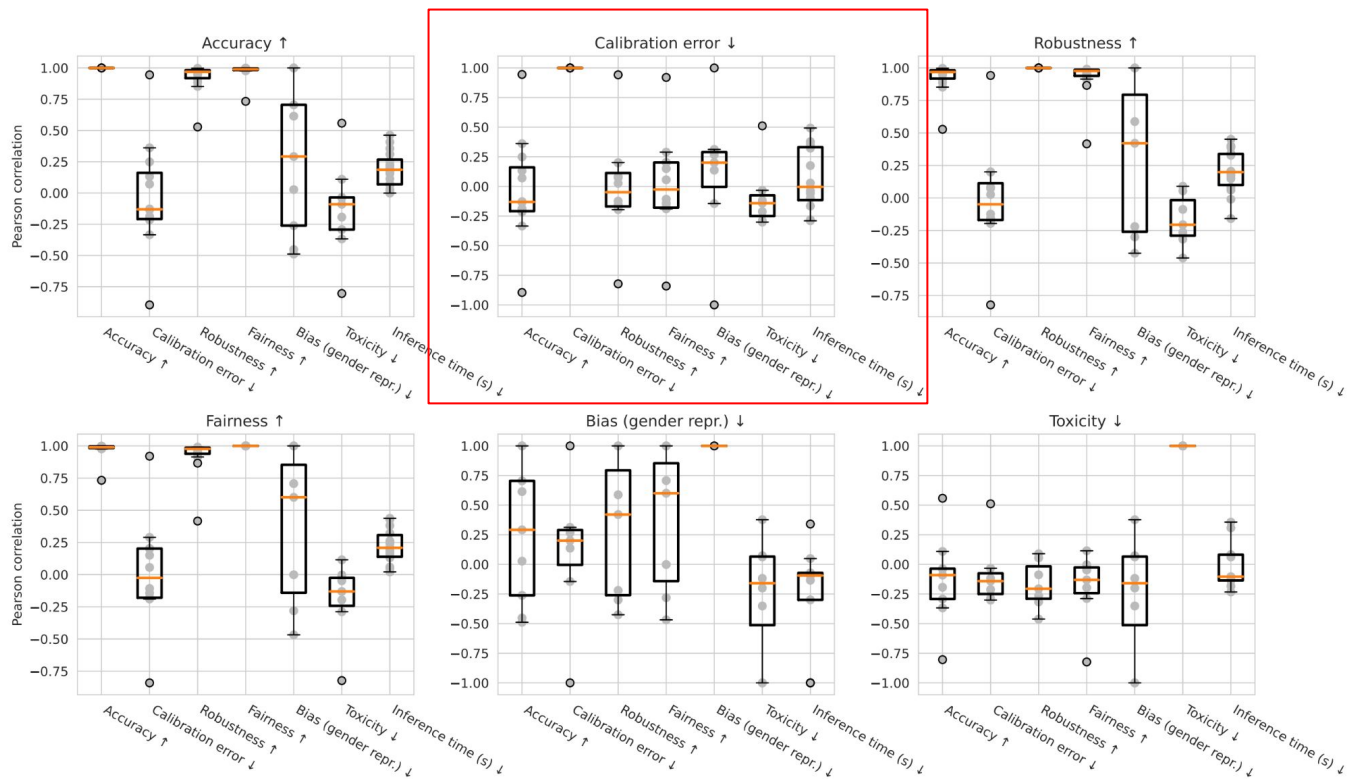
Holistic
Evaluation
of
Language
Models
(HELM)



Correlation between Metrics



Correlation between Metrics



Expected Calibration Error

Confidence

A model's confidence is defined as the $p_\theta(y^*|x)$ where $y^* = \operatorname{argmax}_y p_\theta(y|x)$

ECE

We can obtain ECE by splitting the predictions into bins B_1, \dots, B_M by confidence scores

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{n} |\text{accuracy}(B_m) - \text{confidence}(B_m)|$$

Expected Calibration Error

Probabilities of
model predictions:

0.0

0.1

0.2

0.3



0.7

0.8

0.9

1.0



Equal-sized bins:

Bin 1

Bin 2

$$\text{Accuracy} = 2/4 = 0.5$$

$$\text{Prob} = (0.0 + 0.1 + 0.2 + 0.3) / 4 = 0.15$$

$$\text{Bin-1 error} = |0.5 - 0.15| = 0.35$$

$$\text{Accuracy} = 3/4 = 0.75$$

$$\text{Prob} = (0.7 + 0.8 + 0.9 + 1.0) / 4 = 0.85$$

$$\text{Bin-2 error} = |0.75 - 0.85| = 0.1$$

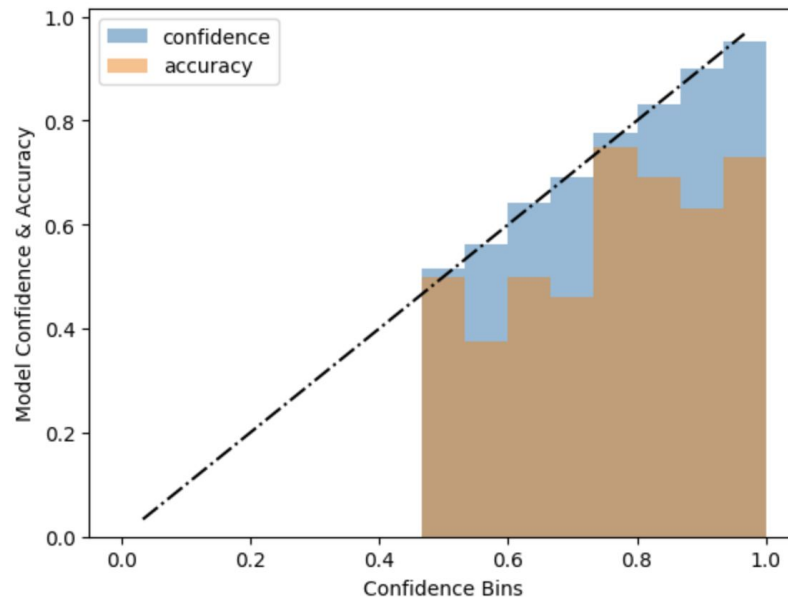
$$\text{ECE (expected calibration error)} = (4/8) * 0.35 + (4/8) * 0.1 = 0.225$$

Expected Calibration Error

Model: BERTforSequence Classification

Task: RTE Dataset in GLUE

Implementation: See Jupyter Notebook



Selective Prediction

Coverage

Given a threshold $t \geq 0$, the coverage $c(t)$ is the fraction of examples for which the model's confidence is at least t

Selective Accuracy

selective accuracy $a(t)$ at a threshold $t \geq 0$ is the accuracy for all examples where the model's predicted confidence is at least t

Selective Prediction

Probabilities of
model predictions:

0.0



0.1



0.2



0.3



0.7



0.8



0.9



1.0



C% (e.g. 10%) of
examples with
highest
probabilities



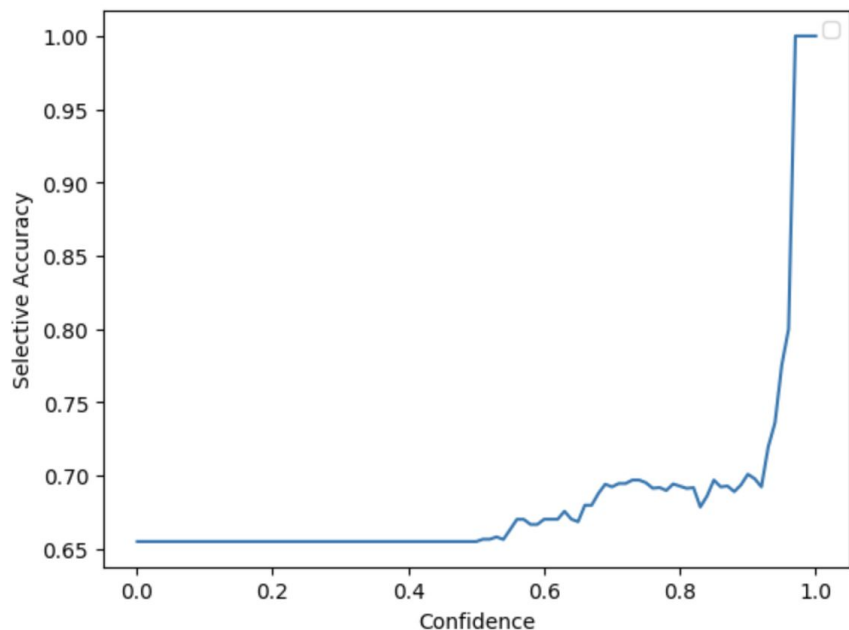
Selective classification accuracy = $2/3 = 0.67$

Selective Coverage-Accuracy Area (SCAA)

Model: BERTforSequence Classification

Task: RTE Dataset in GLUE

Implementation: See Jupyter Notebook



Acknowledgement

This presentation is adapted from Holistic Evaluation of Language Models

(<https://arxiv.org/pdf/2211.09110.pdf>) and Lecture 8

(<https://github.com/nyu-cs2590/course-material/tree/gh-pages/fall2023/lecture/lec08>)