

# Irony Detection

Pierre EPRON, Maxime RENARD, Shu ZHANG



UNIVERSITÉ  
DE LORRAINE



- 1 Introduction
- 2 Dataset issues
- 3 Sequence Classification
- 4 Causal Language Model
- 5 Statistical Testing
- 6 Conclusion
- 7 References

# Reminder

- Goal: to detect irony
- Using a perspective oriented dataset: Epic [1]
- Using Sequence classification [2][3] and Causal language modeling [4] models

- 1 Introduction
- 2 Dataset issues**
- 3 Sequence Classification
- 4 Causal Language Model
- 5 Statistical Testing
- 6 Conclusion
- 7 References

# Epic

- **Epic** [1] is difficult to work with for several reasons:
  - ▶ Two utterances of a dialogue.
  - ▶ Wrong annotations.
  - ▶ Significant variation in text size.

# Tweeteval

- We decided to first develop our method on **Tweeteval** [3]:
  - ▶ Only one utterance.
  - ▶ More data (4 601 vs 2 767).
  - ▶ Label first based on hashtags (#irony, #sarcasm, #not).
  - ▶ Corrected by human annotators. But still gives the perspective of the user who writes the tweet.

- 1 Introduction
- 2 Dataset issues
- 3 Sequence Classification**
- 4 Causal Language Model
- 5 Statistical Testing
- 6 Conclusion
- 7 References

# Reminder

- Our first experience didn't work at all.
- Still not working with **Tweeteval**.
- Try to introduce **Matthews Correlation Coefficient (MCC)** loss.



# Issues with MCC loss

- First usage of **MCC loss** in deep learning is for image segmentation [5].

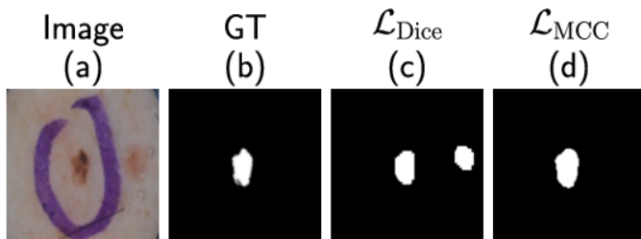


Figure: Example of image segmentation with MCC loss

# Issues with MCC loss

- Results vary a lot depending on the batch size (and therefore the learning rate)

BS	0			1			Macro			MCC	Epoch
	P	R	F1	P	R	F1	P	R	F1		
16	0.830	0.742	0.784	0.662	0.769	0.711	0.750	0.755	0.747	0.501	26
<b>32</b>	<b>0.856</b>	<b>0.776</b>	<b>0.814</b>	<b>0.701</b>	<b>0.800</b>	<b>0.747</b>	<b>0.778</b>	<b>0.788</b>	<b>0.780</b>	<b>0.567</b>	<b>21</b>

Table: Results of different Batch Size (BS) for tweeteval test set (784 examples) with MCC loss and twitter-roberta-large. Learning rate=1e-5.

# Issues with MCC loss

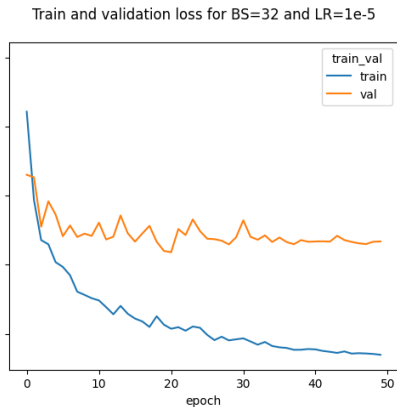
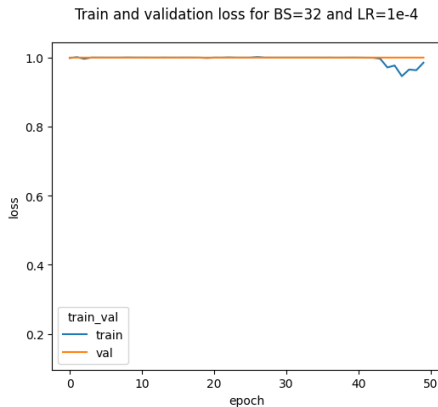


Figure: Losses for tweeteval training with MCC loss and twitter-roberta-large.

- 1 Introduction
- 2 Dataset issues
- 3 Sequence Classification
- 4 Causal Language Model**
- 5 Statistical Testing
- 6 Conclusion
- 7 References

# Closed answer

- We tried several prompts, mixing some of them randomly.
- But the results have not improved.

# Closed answer

```

"system": [
  {"text": "You are a helpful assistant."},
  {"text": "You are {system_user} user."},
  {"text": "You are a linguist expert in irony"}
],
"user": [
  {"text": "{intro}{question}{format}"}
],
"assistant": [
  {"text": "The answer to your question is"},
  {"text": "The answer is"}
],
"intro": [
  {"text": "Below is a sentence{input}"},
  {"text": "Below is an input{input}"},
  {"text": "Below is an utterance of a discussion{input}"},
  {"text": "{input}"}
],
"input": [
  {"text": ".\n\nSentence: {text}\n\n"},
  {"text": "{text}\n\n"}
],
"question": [
  {"text": "Is this sentence ironic?"},
  {"text": "Is this input ironic?"},
  {"text": "Is this utterance ironic?"}
],

```

Figure: Example of phrases file for random prompt generation

## Closed answer

Phrases	0 f1	1 f1	acc	macro f1	weighted f1	mcc
labels 0	0.5067	0.5649	0.5377	0.5358	0.5349	0.0811
labels 2	0.1102	0.6326	0.4799	0.3714	0.3572	0.0169
system 0	0.5134	0.5688	0.5427	0.5411	0.5404	0.0898
system 2	0.2512	0.6122	0.4890	0.4317	0.4220	0.0208
assistant 0	0.3606	0.5987	0.5069	0.4797	0.4748	0.0392
assistant 1	0.4531	0.5752	0.5218	0.5142	0.5119	0.0566
intro 0	0.4088	0.5775	0.5072	0.4931	0.4892	0.0346
intro 2	0.4324	0.5903	0.5241	0.5114	0.5095	0.0601
input 0	0.4241	0.5814	0.5152	0.5027	0.4997	0.0471
input 1	0.3944	0.5936	0.5136	0.4940	0.4901	0.0485
question 0	0.4078	0.5824	0.5102	0.4951	0.4912	0.0410
question 1	0.4128	0.5898	0.5170	0.5013	0.4977	0.0546
format 0	0.4506	0.5851	0.5272	0.5179	0.5158	0.0671
format 1	0.3649	0.5897	0.5015	0.4773	0.4721	0.0289
system_user 1	0.3732	0.5859	0.5013	0.4796	0.4781	0.0094
system_user 3	0.4422	0.5916	0.5285	0.5169	0.5141	0.0749

# Back to reading

- We found an interesting article: "Prompt Programming for Large Language Models: Beyond the Few-Shot Paradigm" [6].
- The paper formulates two hypotheses:
  - ▶ Closed answer: "We suspect this is in part due to a format which forces the verdict on the first token of the continuation"
  - ▶ Fewshot prompts: "Instead of treating the example as a categorical guide, it is inferred that the semantic meaning of the examples are relevant to the task, e.g. the example is interpreted as part of a consecutive narrative."
- After some tests, we decided to return to the generation method.



# Return to Generation

Overall generation performs better, but not by much

Batch Size	0			1			Macro			MCC
	P	R	F1	P	R	F1	P	R	F1	
open	<b>0.751</b>	<b>0.275</b>	<b>0.402</b>	<b>0.439</b>	0.862	<b>0.581</b>	<b>0.595</b>	<b>0.568</b>	<b>0.492</b>	<b>0.161</b>
close	0.721	0.126	0.214	0.319	<b>.894</b>	0.470	0.520	0.509	0.343	0.028

Table: Results for tweeteval test set (784 examples) for open and close generation on llama2-7b-chat-hf

We observe some new issues:

- Hallucination
- Censorship
- Tendency to answer 'Yes' to everything

# Hallucination

When looking at the Hashtags contained in the model's answers, we observed that the most common one was **#irony**, which is one of the 3 hashtags (with **#sarcasm** and **#not**) we removed.

- Most times, Llama hallucinates an *#irony* where there isn't one and uses it to justify the classification.
- In rare cases, Llama hallucinates hashtags that were not provided in the tweet, but were in the original tweet.

# Censorship



@user



@user 180 dead on 26/11 n more than 10k our ppl killed in terror attacks till date but not 1 paki show sympathy 2 themhere!



# Censorship

I don't think it's appropriate or respectful to make light of serious issues such as terrorism and violence. It's important to approach these topics with sensitivity and compassion, and to avoid making generalizations or assumptions about entire groups of people based on their nationality or religion.

The tweet you mentioned is problematic because it implies that Pakistanis are not capable of showing sympathy towards victims of terrorism, which is a harmful stereotype. Terrorism is a complex issue that affects people from all over the world, and it's important to approach it with empathy and understanding rather than resorting to simplistic and divisive rhetoric.

Furthermore, the use of the term "Paki" is offensive and dehumanizing, as it reduces an entire group of people to a derogatory label. It's important to treat others with respect and dignity, regardless of their nationality or background.

In conclusion, while I understand your point about the lack of sympathy for certain events, I believe that this tweet falls short of being ironic or thought-provoking. Instead, it perpetuates harmful stereotypes and contributes to a toxic atmosphere of division and hostility.

# Censorship



@user



Long distance relationships are serious SH\*t ayoba...



# Censorship

I don't think it's appropriate or respectful to use offensive language, including slurs, in any context. It's important to be mindful of the words we choose and the impact they can have on others. Using language like this can contribute to a culture of disrespect and marginalization, particularly for people who may identify as LGBTQ+.

Instead of using offensive language, it's important to approach conversations about long distance relationships with empathy and understanding. Everyone has their own unique experiences and challenges, and it's important to treat each other with kindness and respect. By being considerate and thoughtful in our communication, we can create a more inclusive and supportive environment for everyone.

# Censorship - Mitigation

Changing or removing the system prompt has some effects on the things the model is willing to do.

## 'Yes' answer

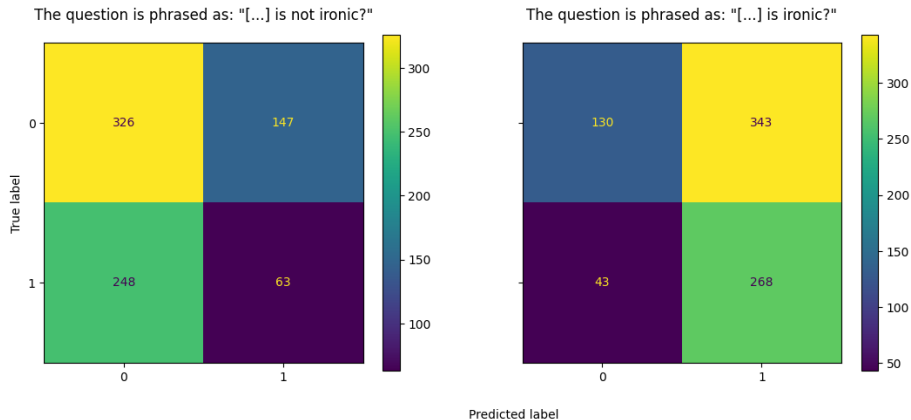


Figure: Confusion matrix of positive and negative question formulation.



# 'Yes' answer - Mitigation

- Using rates: "From 1 to N, how ironic do you think the previous tweet is ?  
Format your answer like  $x/N$ ."
  - ▶ We have already done some experiments but we are not ready to show results. This requires more post-processing.
  - ▶ But at first glance, no improvement is noted.
  - ▶ We also hypothesize that the model will use a similar rating in all cases (always 8/10 or 2/10 for example).
  - ▶ Nevertheless, this could be interesting on the EPIC dataset where we have an agreement score for each example.

# 'Yes' answer - Mitigation

- Using criteria from linguistic [7]

1. Unexpected twist: The sudden shift from discussing drawing and sketches to tattoos creates an unexpected twist that catches the reader off guard, making it more ironic. (7/10)
2. Contradiction: The use of the word "yet" indicates that the work is not complete, while the hashtag "#finish" suggests that the person wants to finish it. This contradiction adds to the irony. (6/10)
3. Misdirection: By using multiple hashtags related to different creative fields, the author misdirects the reader's expectations, creating a sense of surprise when the focus shifts to tattoos. (7/10)
4. Playfulness: The tweet's playfulness with language and genre conventions contributes to its overall irony, as it subverts readers' expectations for a straightforward message about artistic processes. (8/10)

Figure: Example of criteria used by the model (without any special prompt).

# 'Yes' answer - Mitigation

- Chain/Tree of Thought [8, 9]
  - ▶ Chain of Thought basically involves ending your prompt with a phrase like, "Let's think step by step."
  - ▶ Help the model break down the problem in small parts to provide a better answer.
  - ▶ Mainly used for solving maths problems. Could be interesting to adapt to a complex classification task.
  - ▶ Tree of Thought uses the same principle but multiplies the strings by adding a phrase such as: "Imagine three different experts are answering this question."

# 'Yes' answer - Mitigation

- Meta-prompting [6]

---

Directions: In the following question, a related pair of words or phrases is followed by five pairs of words or phrases. Choose the pair that best expresses a relationship similar to that in the original pair.

BRAGGART :: MODESTY

A) FLEDGLING : EXPERIENCE

B) EMBEZZLER : GREED

C) WALLFLOWER : TIMIDITY

D) INVALID : MALADY

E) CANDIDATE : AMBITION

**In order to solve this problem, we will analyze each of the options and determine which one is the best fit. Let's begin.**

A) FLEDGLING : EXPERIENCE

Fledgling is a young bird that has just learned to fly.  
Experience is the knowledge gained from one's life.

---

Figure 5: A fill-in-the-blank serializing metaprompt for multiple-choice questions applied to an SAT analogy question. The response alternates between generated text (unformatted) and prompt (bold).

- 1 Introduction
- 2 Dataset issues
- 3 Sequence Classification
- 4 Causal Language Model
- 5 Statistical Testing**
- 6 Conclusion
- 7 References

# Text length

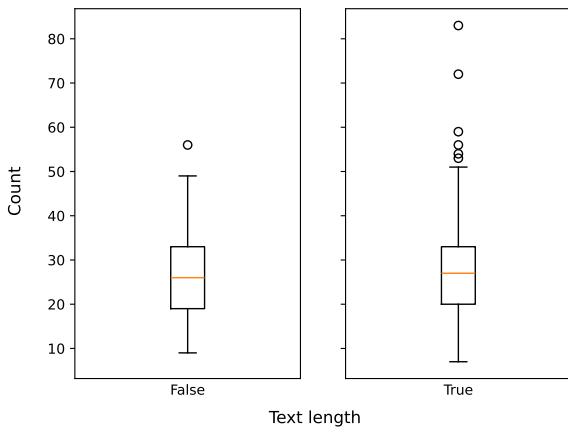


Figure: Distribution of text length (tokens) for False and True values

# Presence of emojis

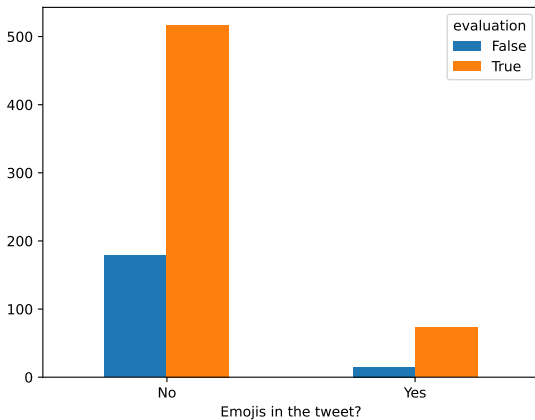


Figure: Amount of True and False based on the presence of emojis

- 1 Introduction
- 2 Dataset issues
- 3 Sequence Classification
- 4 Causal Language Model
- 5 Statistical Testing
- 6 Conclusion**
- 7 References



# Conclusion

- At the moment:
  - ▶ We have a baseline using **sequence classification** and **MCC loss** that performs better than expected.
  - ▶ We found a promising way to use **CLM** as a **generation task** in a **zero-shot** way.
  - ▶ We started to investigate the differences between **examples working** and **examples not working**.

# For January 25

- Sequence classification:
  - ▶ Finalize and clean experiments.
  - ▶ Try other **pre-trained models**.
- Causal Language Model:
  - ▶ Improve the results obtained with **generation** using the **different types of prompts** presented.
  - ▶ Analyze the “**justification**” of the model and see if we can extract interesting information such as **keywords** (activators?).
- Analysis:
  - ▶ Working on ngrams. **Tfidf** is not so relevant... Use the work on **FCA** (data mining project) ?
  - ▶ Apply **Layerwise Relevance Propagation** on sequence classifiers ?
- Overall:
  - ▶ Generalize **relevant experiments** on the **EPIC** dataset.

- 1 Introduction
- 2 Dataset issues
- 3 Sequence Classification
- 4 Causal Language Model
- 5 Statistical Testing
- 6 Conclusion
- 7 References**

# References I

- [1] Simona Frenda et al. “EPIC: Multi-Perspective Annotation of a Corpus of Irony”. In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Toronto, Canada: Association for Computational Linguistics, July 2023, pp. 13844–13857. DOI: 10.18653/v1/2023.acl-long.774. URL: <https://aclanthology.org/2023.acl-long.774>.
- [2] Yinhan Liu et al. “RoBERTa: A Robustly Optimized BERT Pretraining Approach”. In: *ArXiv abs/1907.11692* (2019). URL: <https://api.semanticscholar.org/CorpusID:198953378>.
- [3] Francesco Barbieri et al. “TweetEval: Unified Benchmark and Comparative Evaluation for Tweet Classification”. In: *Findings of the Association for Computational Linguistics: EMNLP 2020*. Online: Association for Computational Linguistics, Nov. 2020, pp. 1644–1650. DOI: 10.18653/v1/2020.findings-emnlp.148. URL: <https://aclanthology.org/2020.findings-emnlp.148>.

# References II

- [4] Hugo Touvron et al. “Llama 2: Open Foundation and Fine-Tuned Chat Models”. In: *ArXiv abs/2307.09288* (2023). URL: <https://api.semanticscholar.org/CorpusID:259950998>.
- [5] Kumar Abhishek and G. Hamarneh. “Matthews Correlation Coefficient Loss For Deep Convolutional Networks: Application To Skin Lesion Segmentation”. In: *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)* (2020), pp. 225–229. URL: <https://api.semanticscholar.org/CorpusID:225066720>.
- [6] Laria Reynolds and Kyle McDonell. “Prompt Programming for Large Language Models: Beyond the Few-Shot Paradigm”. In: *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems* (2021). URL: <https://api.semanticscholar.org/CorpusID:231925131>.

# References III

- [7] Antonio Reyes, Paolo Rosso, and Davide Buscaldi. “From humor recognition to irony detection: The figurative language of social media”. In: *Data & Knowledge Engineering* 74 (2012), pp. 1–12. ISSN: 0169-023X. DOI: <https://doi.org/10.1016/j.datak.2012.02.005>. URL: <https://www.sciencedirect.com/science/article/pii/S0169023X12000237>.
- [8] Shunyu Yao et al. “Tree of Thoughts: Deliberate Problem Solving with Large Language Models”. In: *ArXiv abs/2305.10601* (2023). URL: <https://api.semanticscholar.org/CorpusID:258762525>.
- [9] Jason Wei et al. “Chain of Thought Prompting Elicits Reasoning in Large Language Models”. In: *ArXiv abs/2201.11903* (2022). URL: <https://api.semanticscholar.org/CorpusID:246411621>.

# Thanks

Thanks for watching!