

Using BERT to investigate cognitive biases in thinking of everyday unexpected events

1 | Example Language Model Use Case

**2 | Resources for Getting a Jump-Start on Text
Classification with BERT**

Contents

1 | Predicting Unexpected Events

2 | Classifying Unexpected Events

3 | Explaining the Classifications

4 | Resources

5 | Conclusion

Predicting Unexpected Events

In a foreign city and worried about missing a flight?

Q: What unexpected things could occur?

- Traffic is really bad.
- Buses are unsafe.
- Taxis are unreliable.
- Police outriders escort me to the airport.

What is the Unexpected?

- Subjective probability?
- What objects do people mutate in the unexpected?
- What are the features of the event?
Valence, controllability, goals...?

Example Material

Sentence Type	Scenario
Goal	Louise wants to shop at an expensive clothes store.
Intermediate Event	She is wearing her favourite dress and matching shoes.
Plan Step	Louise draws money from the ATM.

Instructions

Unexpected: Something unexpected occurred. What do you think happened?

Expected: Assuming nothing unexpected occurs, what do you expect to typically happen next?

Example Responses

Unexpected

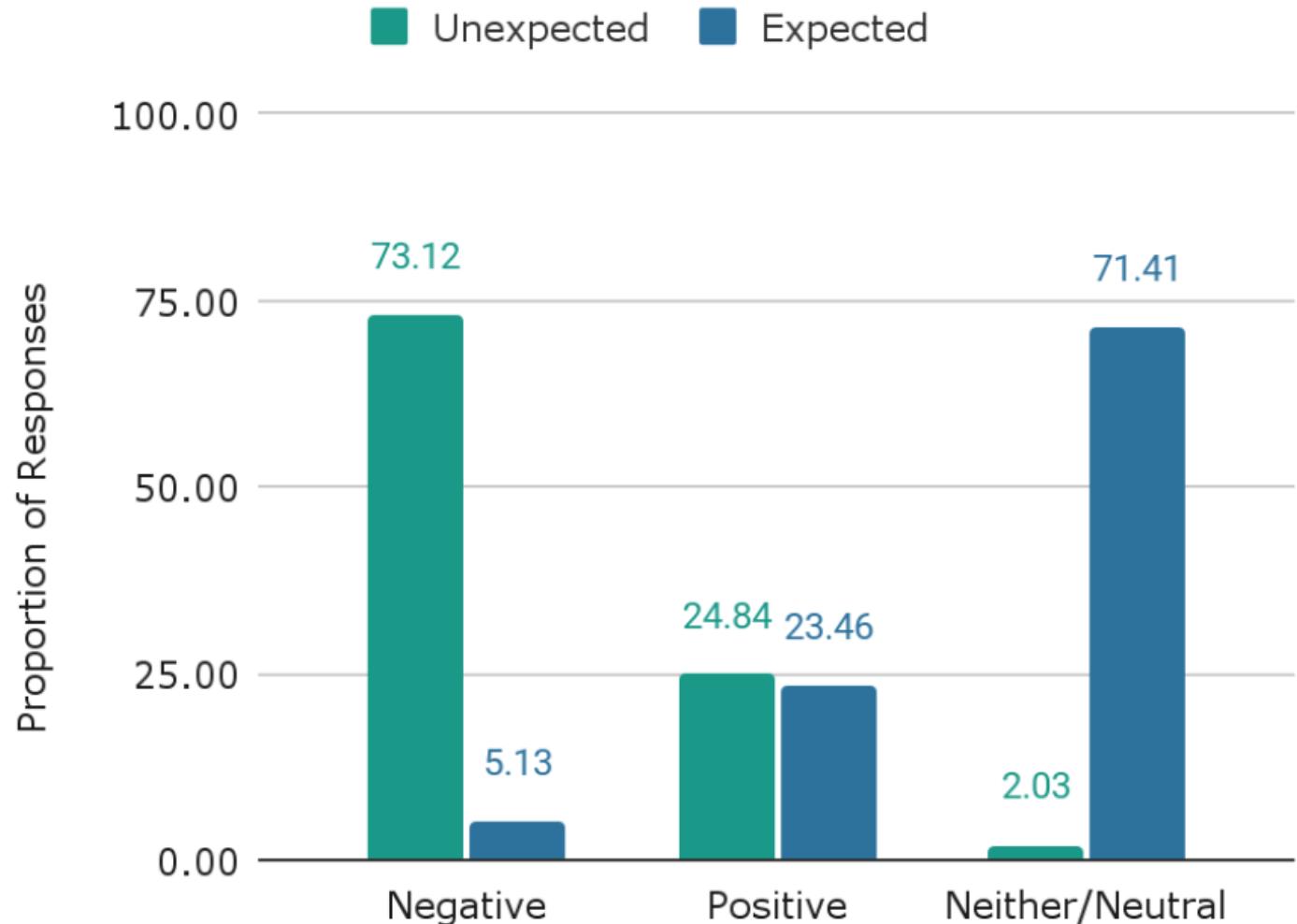
- “She saw the account balance as she withdrew money and realized she didn't have enough to justify the shopping trip”
- “She is robbed at the ATM”
- “She had more money than she realised so could buy something expensive.”

Expected

- “She will go into the expensive clothes store and look around. If there is something she likes she may try it on and she may buy some clothes”
- “I would expect Louise to take her money to the clothes store and spend her money.”
- “she buys something to match her outfit”

Valence

Prediction of unexpected events is valenced and skewed towards the **negative**.



Contents

1 | Predicting Unexpected Events

2 | Classifying Unexpected Events

3 | Explaining the Classifications

4 | Resources

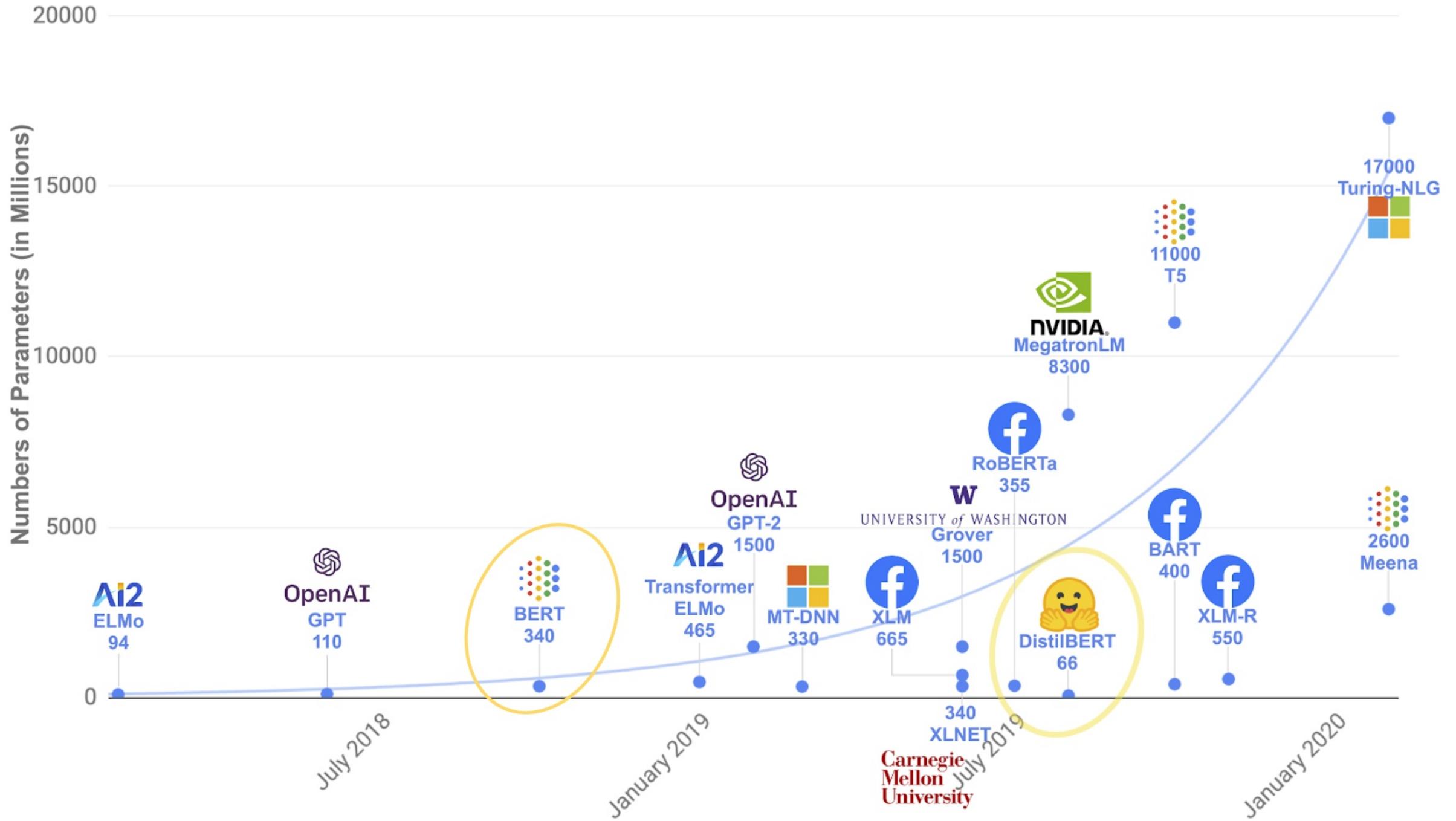
5 | Conclusion

What is the Goal?

- Some say that language models can serve as models of “psychological semantics”
 - Meaning computational word representations are similar to human mental word representations
- But others say this is the wrong way to look at it.

“... our critique is that the current [Language Models] are too strongly linked to complex text-based patterns, and too weakly linked to world knowledge”

LAKE, B. M. & MURPHY, G. L. (2021). WORD MEANING IN MINDS AND MACHINES. PSYCHOLOGICAL REVIEW. ADVANCE ONLINE PUBLICATION.
[HTTPS://DOI.ORG/10.1037/REV0000297](https://doi.org/10.1037/REV0000297)

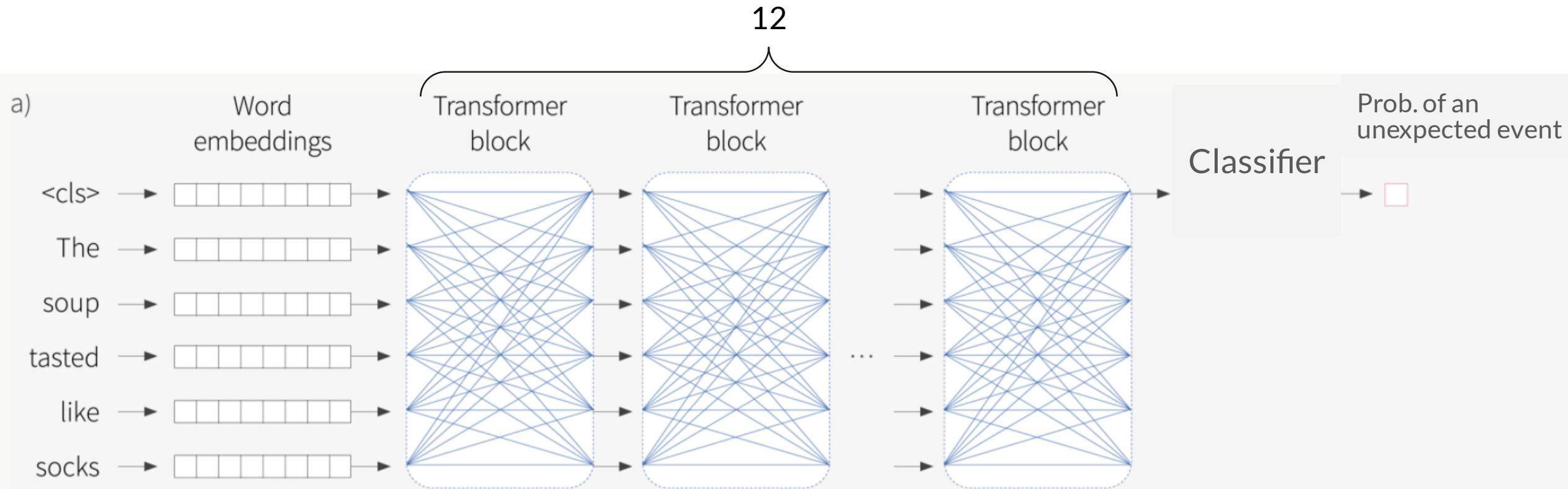


Why BERT, Specifically?

- Short texts (max 512 tokens)
- Convenient! (open-source)

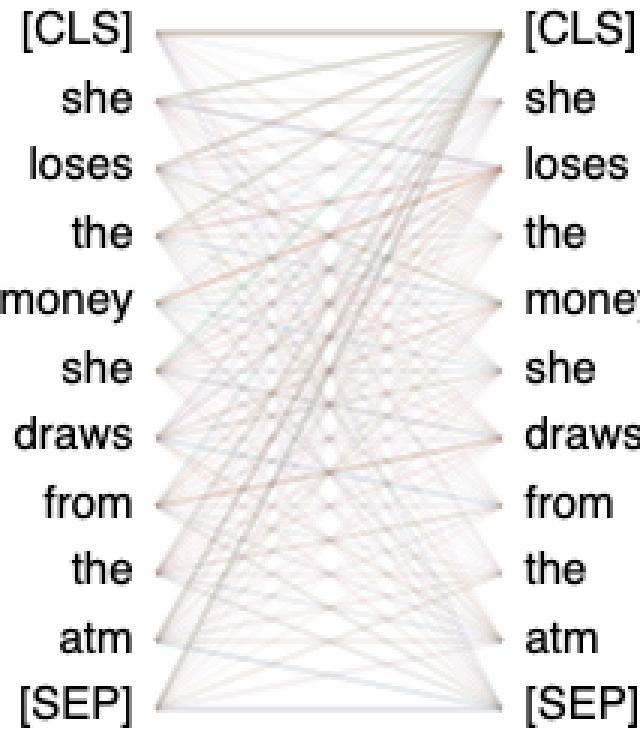
How Does BERT Work? a cursory runthrough

BERT Layers

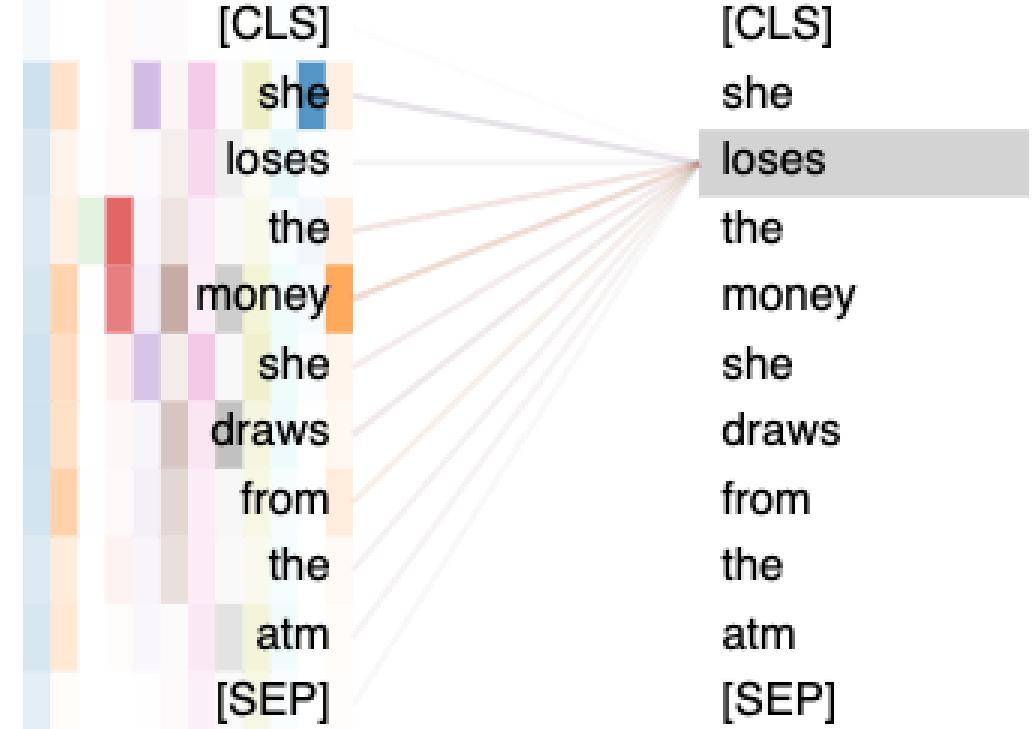


Each word is informed, to some degree, by every word

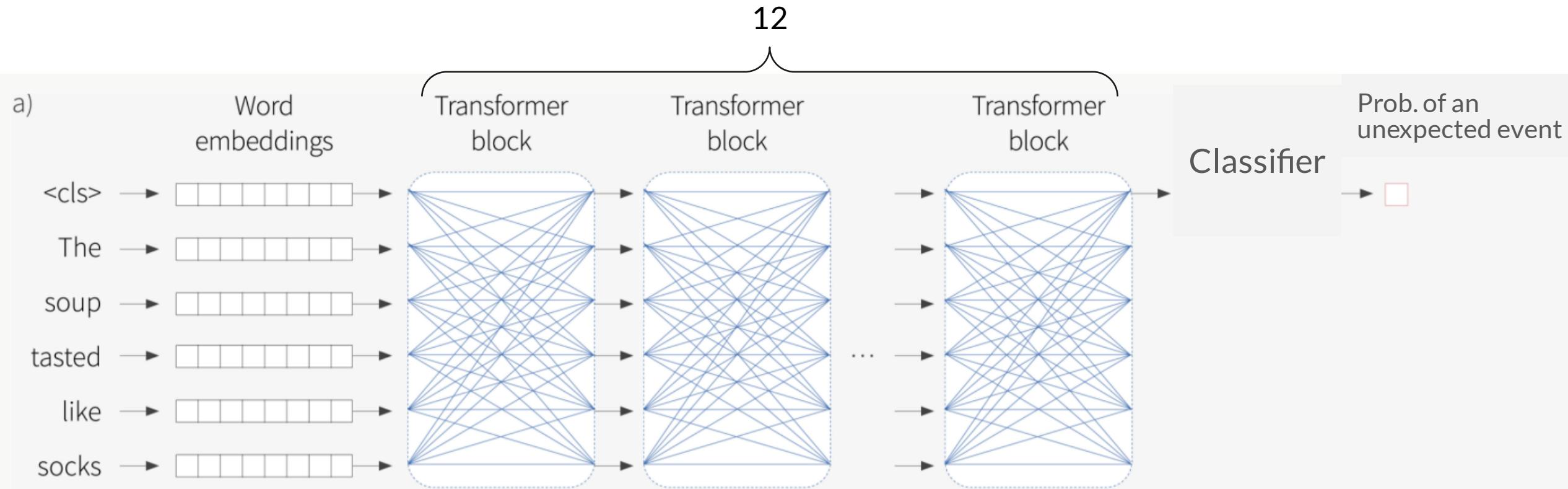
All Attention Weights



Attention Weights for “loses”



BERT Layers



Resources



huggingface Transformers

<https://huggingface.co/docs/transformers/>



BERT for Humanists

<http://www.bertforhumanists.org/>

<https://www.youtube.com/watch?v=UmyOhI9Acil>

[https://colab.research.google.com/drive/19jDqa5D5XfxPU6NQef17BC07xQdRnaKU?
usp=sharing](https://colab.research.google.com/drive/19jDqa5D5XfxPU6NQef17BC07xQdRnaKU?usp=sharing)

Current Experiments

- 1 | Logistic Regression (baseline)
- 2 | Fine-Tuned BERT
- 3 | Fine-Tuned DistilBERT
- 4 | LIME Explanations

Train/Test Split

Test Set	Training Subset 1	Training Subset 2	Training Subset 3
steve_gardening	rebecca_swimming	katie_kitten	anna_interview
louise_shopping	sally_wine	lucy_loan	michael_tea
alan_plane	karen_bus	belinda_meeting	robert_essay
edith_exam	bob_job	peter_college	sam_driving
mary_food	bill_holiday	john_party	sean_call

General Method

1 | Logistic Regression

- Using BERT (WordPiece) Encodings
- TF-IDF Vectors
- Binary logistic regression model
- Trained on all 3 training sets
- Tested on test set

2 | BERT Models

- Hyperparameter search by crossvalidation:
- Fine-Tune on 2 training sets, and validate on the 3rd training set
- Choose best hyperparameters for final model training
- Final Models:
- Fine-Tuned on all 3 training sets
- Tested on test set

Results

	Logistic Regression			BERT Not Fine-Tuned			Support
	Precision	Recall	F1-Score	Precision	Recall	F1-Score	
expected	0.89	0.68	0.77	0	0	0	195
unexpected	0.70	0.89	0.78	0.45	1.00	0.62	160
accuracy			0.78			0.45	355
macro avg	0.79	0.79	0.78	0.23	0.50	0.31	355
weighted avg	0.80	0.78	0.78	0.20	0.45	0.28	355
	DistilBERT Fine-Tuned			BERT Fine-Tuned			Support
	Precision	Recall	F1-Score	Precision	Recall	F1-Score	
expected	0.95	0.85	0.90	0.92	0.87	0.89	195
unexpected	0.84	0.95	0.89	0.85	0.91	0.88	160
accuracy			0.90			0.89	355
macro avg	0.90	0.90	0.90	0.89	0.89	0.89	355
weighted avg	0.90	0.90	0.90	0.89	0.89	0.89	355

Contents

1 | Predicting Unexpected Events

2 | Classifying Unexpected Events

3 | Explaining the Classifications

4 | Resources

Logistic Regression: TF-IDF Weighted Coefficients

Unexpected

Expected

When he arrives from his trip, he starts to go through his luggage to find his new shirt. He finds it, but it's not the right size. He goes to the store to buy a new one. While he's there, he sees a sale on shirts, so he buys two more. After buying the shirts, he goes to eat at a restaurant with his friends. They order food and have a great time. When they're finished eating, they go to a bar to drink some beer. They have fun at the bar, and then they go home. The next day, he wakes up and gets ready for work. He takes a shower and gets dressed. He leaves for work and has a good day at the office. At the end of the day, he goes home and relaxes. He feels tired but happy.

LIME (Locally Interpretable Model-Agnostic Explanations)

id	text	
1 (input)	Someone robs her	
2	[MASK] robs her	
3	Someone [MASK] her	
...	...	
n	Someone robs [MASK]	

LIME (Locally Interpretable Model-Agnostic Explanations)

id	text	prediction
1 (input)	Someone robs her	unexpected
2	[MASK] robs her	unexpected
3	Someone [MASK] her	expected
...
n	Someone robs [MASK]	unexpected

LIME (Locally Interpretable Model-Agnostic Explanations)

id	text	prediction
1 (input)	Someone robs her	unexpected
2	[MASK] robs her	unexpected
3	Someone [MASK] her	expected
...
n	Someone robs [MASK]	unexpected

BERT Average LIME Weights

Unexpected

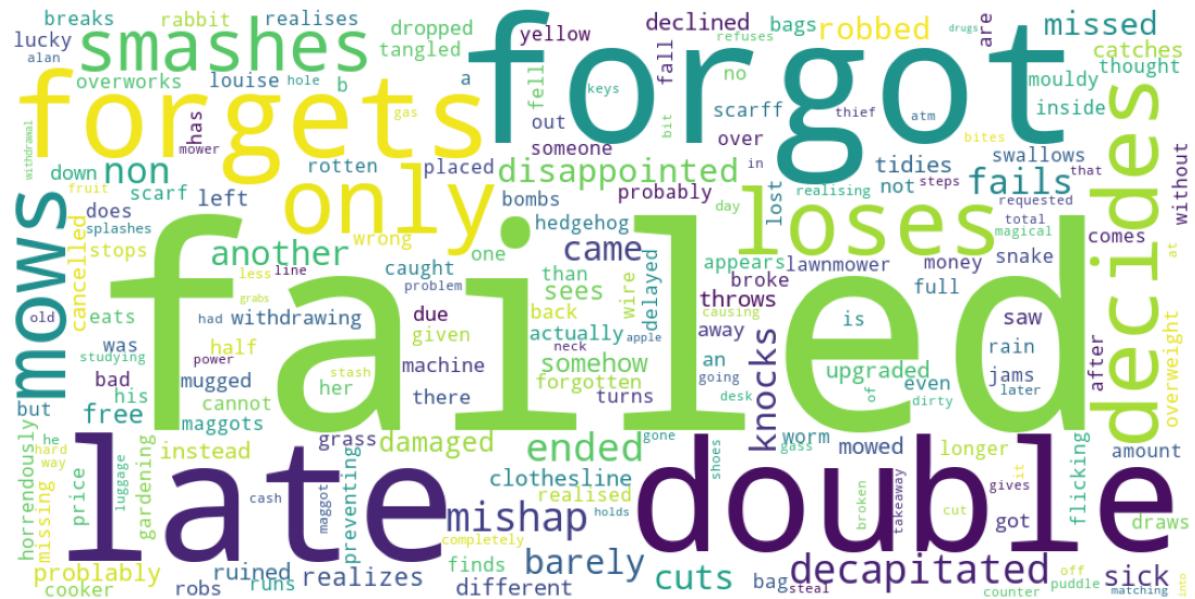


Expected



DistilBERT Average LIME Weights

Unexpected



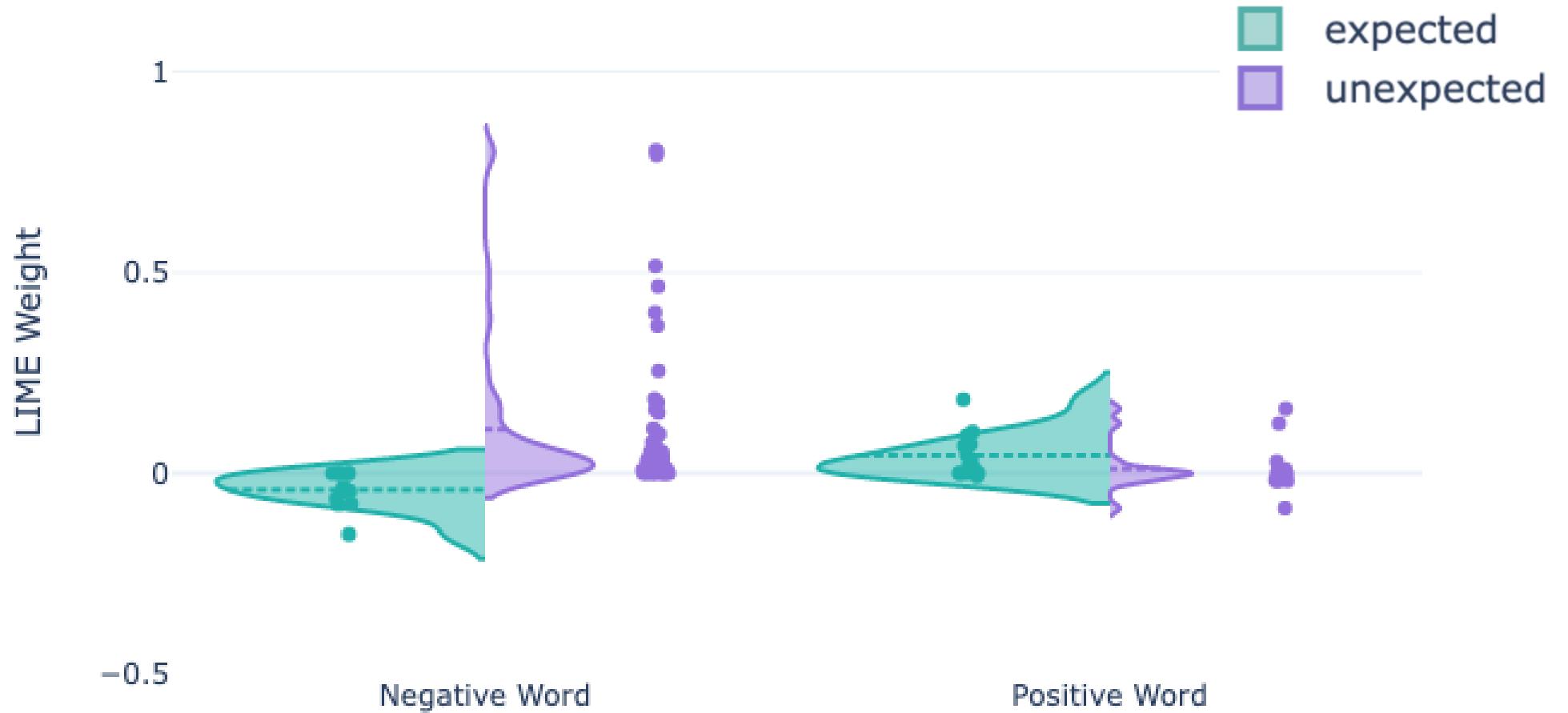
Expected

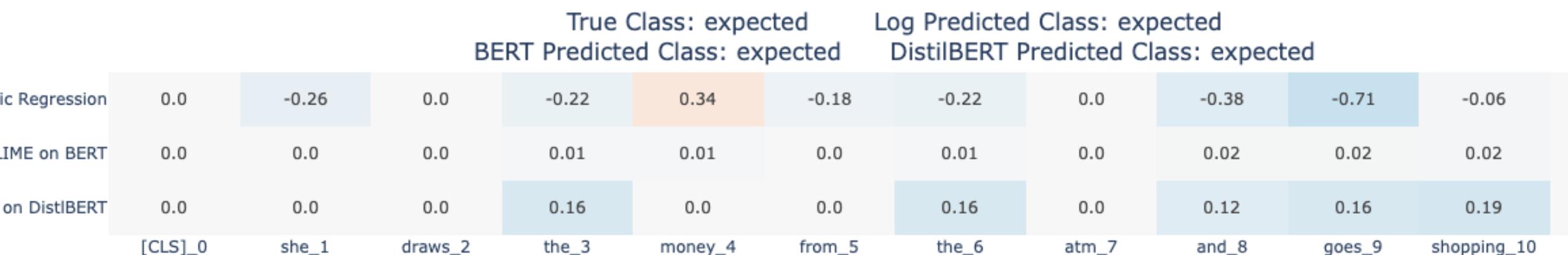


BERT Classifications



DistilBERT Classifications





Contents

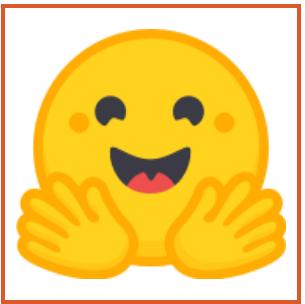
1 | Predicting Unexpected Events

2 | Classifying Unexpected Events

3 | Explaining the Classifications

4 | Resources

Resources



huggingface Transformers

<https://huggingface.co/docs/transformers/>



BERT for Humanists

<http://www.bertforhumanists.org/>

[https://www.youtube.com/watch?
v=UmyOhI9Acil](https://www.youtube.com/watch?v=UmyOhI9Acil)

[https://colab.research.google.com/drive/
19jDqa5D5XfxPU6NQef17BC07xQdRnaKU?
usp=sharing](https://colab.research.google.com/drive/19jDqa5D5XfxPU6NQef17BC07xQdRnaKU?usp=sharing)



Interpretable Machine Learning

[https://christophm.github.io/
interpretable-ml-book/](https://christophm.github.io/interpretable-ml-book/)

Thank you!

- 1 | People believe the unexpected is negative and uncontrollable.
- 2 | BERT models this negativity bias during fine-tuning.

Contact me:



molly.quinn@ucdconnect.ie



@Molly_S_Quinn



scholar.google.com/citations?user=M677etQAAAAJ

Watch Dr. Maria Antoniak's Tutorial:

BERT for Computational Social Scientists

<https://www.youtube.com/watch?v=UmyOhl9Acil>

Limitations

- May be overfitting to training data
- Explanation method does not explain the model itself
- Questionable generalizability
 - to other datasets
 - other languages, and
 - behavioral/cultural factors

Future Directions

- Other explanation methods...
- Generating unexpected events
- Explicitly incorporating sentiment/etc. into classification/generation of unexpected events.