

Experiment Protocol for CS224U, Spring 2020.

Abhishek Goswami
Microsoft / Redmond, WA
agoswami@microsoft.com

Abstract

This document sets the experiment protocol for CS224U, Spring 2020 quarter.

1 Hypothesis

Natural language inference (NLI) is the task of determining if a natural language hypothesis can be inferred from a given premise in a justifiable manner. Existing models perform well at standard evaluations for NLI, achieving impressive results in leaderboards such as GLUE. However, a growing body of evidence (McCoy et al., 2019; Glockner et al., 2018) shows that state-of-the-art models learn to exploit spurious statistical patterns in datasets instead of learning meaning in the flexible and generalizable way. NLI adversarial evaluations, where existing state-of-the-art NLI systems are evaluated on a completely unseen new test dataset, further exposes these concerns. State-of-the-art NLI systems perform quite poorly in the adversarial evaluation setting (Nie et al., 2019), reflecting that such systems do not represent true competence at natural language understanding. This begs the question: If a system fails an NLI adversarial evaluation, is it a failing of the model or of the dataset used to develop the model?

In this work, we hypothesize that the failings in the adversarial evaluation setting come from the dataset used to develop the model. Stan-

dard evaluations in NLI typically use datasets generated using a single generation process. We argue that these standalone, per-dataset generation processes encode latent signals in the premise/hypothesis construction. This leads to favorable results in the standard evaluation setting, but fail spectacularly in the adversarial evaluation setting. We believe this can be improved by using a mix of multiple datasets during the training process. This setting enables us to leverage several related sentence pair datasets, while avoiding per-dataset statistical concerns.

2 Data

In this section we provide a description of the datasets that we use in this project. We also describe also how we cast some of these datasets for the NLI task, as a form of weak supervision. Table 1 gives a summary of the distribution of labels across different datasets.

2.1 Test Set

For our test set we use the Adversarial NLI dataset, as described in (Nie et al., 2019). More specifically we choose the test set from Round 1. By construction, this test set poses a serious challenge to our models. The examples in this dataset were verified as correct by human annotators. But a state-of-the-art BERT_{LARGE} model trained on SNLI (Bowman et al., 2015)

Categories	Training Datasets							Dev Set	Test Set
	MNLI	MRPC	QQP	STSB	QNLI	RTE	WNLI	MNLI-(m/mm)	ANLI-A1
Contradiction	131k	1194	229k	1773	52k	1241	312	3,213 / 3,240	333
Entailment	131k	2474	134k	1052	52k	1249	323	3,479 / 3,463	334
Neutral	131k	-	-	2924	-	-	-	3,123 / 3,129	333

Table 1: Distribution of labels across different datasets

and MNLI (Williams et al., 2017) datasets got them all wrong. Also, in the spirit of adversarial evaluation, we do not consider any other ANLI datasets during training.

2.2 Dev Set

For our dev set we use both the matched and mismatched dev sets which are part of the MNLI dataset. We chose this as our dev set because it contains all the three different labels of interest. Please do note that since we are interested in the adversarial evaluation setting, the choice of dev set is purely arbitrary.

2.3 Train Set

We give a brief summary of the different datasets that we consider in our mix for training. The common theme for all these datasets is that they lend themselves to pairwise text classification. We argue that each of them can be cast for the NLI setting.

MNLI: An indoor path-loss propagation model essentially forms the bedrock for these techniques.

MRPC: An indoor path-loss propagation model essentially forms the bedrock for these techniques.

QQP: An indoor path-loss propagation model essentially forms the bedrock for these techniques.

STSB: An indoor path-loss propagation model essentially forms the bedrock for these

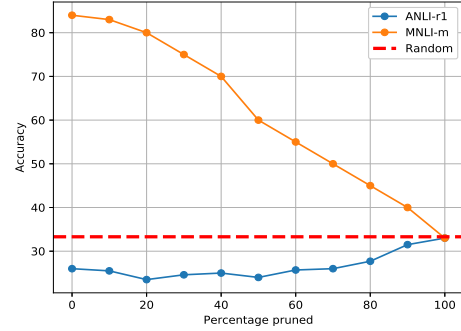


Figure 1: Accuracy as a function of attention heads. BERT_{BASE} model trained on MNLI only.

techniques.

QNLI: An indoor path-loss propagation model essentially forms the bedrock for these techniques.

RTE: An indoor path-loss propagation model essentially forms the bedrock for these techniques.

WNLI: An indoor path-loss propagation model essentially forms the bedrock for these techniques.

3 Metrics

We measure our model performance using the accuracy metric. As shown in Table 1, our dev and test sets have a uniform distribution across all the three labels. We believe this justifies using accuracy as the evaluation metric.

System	Training Data Ratio							Dev Acc MNLI-(m/mm)	Test Acc ANLI-A1
	MNLI	MRPC	QQP	STSB	QNLI	RTE	WNLI		
Random	1.0	-	-	-	-	-	-	33.8 / 33.2	33.3
Baseline	1.0	-	-	-	-	-	-	50.3 / 49.6	28.9
BERT _{BASE}	1.0	-	-	-	-	-	-	84.1 / 84.4	26.0
RoBERTa _{BASE}	1.0	-	-	-	-	-	-	87.7 / 87.6	28.5
RoBERTa _{LARGE}	1.0	-	-	-	-	-	-	92.0 / 90.0	46.0
BERT _{BASE} with weak supervision	0.1	0.1	0.1	-	-	-	-	76.9 / 77.2	23.2
	0.1	0.1	0.1	0.1	0.1	0.1	0.1	77.3 / 77.6	25.8
RoBERTa _{BASE} with weak supervision	0.1	0.1	0.1	-	-	-	-	83.6 / 84.1	26.0
	0.1	0.1	0.1	0.1	0.1	0.1	0.1	83.6 / 84.3	29.3

Table 2: Evaluation results on the Dev/Test sets

Categories	BERT _{BASE}	RoBERTa _{BASE}
Contradiction	369	353
Entailment	303	308
Neutral	328	339

Table 3: Distribution of predicted labels on ANLI-A1. Models trained on MNLI only.

4 Models

A description of the models that you’ll be using as baselines, and a preliminary description of the model or models that will be the focus of your investigation. At this early stage, some aspects of these models might not yet be worked out, so preliminary descriptions are fine.

5 General Reasoning

An explanation of how the data and models come together to inform your core hypothesis or hypotheses.

6 Progress Summary

In this section, we report the progress we have made so far. Table 2 shows a comprehensive view of the experiments done so far.

6.1 Training on MNLI only

A random guessing strategy gives a accuracy of 33.3% on both the Dev and Test sets which

is intuitive given that we have 3 labels which are uniformly distributed (See Section 2).

Compared to Random, our Baseline model does better on the Dev set indicating it has picked up signal from the MNLI training data. However, it does worse than Random in the adversarial evaluation setting when evaluated on the ANLI-A1 test set. We see similar behaviour for both the ‘simple’ Transformer-based models, . When trained on MNLI data both BERT_{BASE} and RoBERTa_{BASE} do remarkably well on the Dev set, but do quite poorly on the ANLI-A1 test set. In this setting, only RoBERTa_{LARGE} is able to outperform Random on the Test set.

6.2 Analyzing Transformer models on Adversarial data

We try to understand what Transformed-based models are doing from two perspectives (a) Error analysis and (b) Role of the attention heads. For this we continue to use models trained on MNLI only.

Error Analysis: We look at some examples where *both* BERT_{BASE} and RoBERTa_{BASE} are wrong, while being in perfect agreement with each other. Table 4 shows one example from each such category. This leads to a few observations:

Premise	Hypothesis	gold_label	BERT _{BASE}	RoBERTa _{BASE}
Julian Peter McDonald Clary (born 25 May 1959) is an English comedian and novelist. Openly gay, Clary began appearing on television in the mid-1980s and became known for his deliberately stereotypical camp style. Since then he has also acted in films, television and stage productions, and was the winner of <i>Celebrity Big Brother</i> 10 in 2012.	Born on May 25, 1958 Julian Clary knew he would be gay.	contradiction	entailment	entailment
Bassingham is a village and civil parish in the North Kesteven district of Lincolnshire, England. The population of the civil parish at the 2011 census was 1,425. The village is situated approximately 8 mi south-west from the city and county town of Lincoln.	Bassingham is a village and civil parish in the North Kesteven district of Lincolnshire, England and only has 1,435 citizens as of the 2011 census.	contradiction	neutral	neutral
Idrees Kenyatta Walker (born February 1, 1979) is a former professional American football player who was an offensive tackle in the National Football League (NFL) for six seasons. Walker played college football for the University of Florida. A first-round pick in the 2001 NFL Draft, he played professionally for the Tampa Bay Buccaneers of the NFL.	Kenyatta Walker did not play football at Florida State.	entailment	contradiction	contradiction
The 2005 Big East Men's Basketball Championship was played from March 9 to March 12, 2005. The tournament took place at Madison Square Garden in New York City. The Syracuse Orange won the tournament and were awarded an automatic bid to the 2005 NCAA Men's Division I Basketball Tournament.	The tournament took place over four days.	entailment	neutral	neutral
Puss in Boots is an action game based on the DreamWorks Animation SKG movie of the same name. It was developed by Blitz Games, and released by THQ on October 25, 2011 for Xbox 360, PlayStation 3, Wii and Nintendo DS. It features support for Kinect and PlayStation Move on the respective platforms. It was released on October 25, 2011 in North America and December 2 for Europe.	Dreamworks released a movie in 2012	neutral	contradiction	contradiction
Binani Industries Ltd is an Indian business group based in Mumbai. It is a 143-year old business conglomerate and belongs to the legendary Braj Binani Group. The business portfolio of Binani Industries includes sectors like cement, zinc, glass-fiber, and downstream composite products.	Braj Binani Group has enjoyed ownership of Binani Industries Ltd since its creation.	neutral	entailment	entailment

Table 4: Examples where *both* BERT_{BASE} and RoBERTa_{BASE} are wrong, while being in perfect agreement with each other. Models trained on MNLI only.

- Adversaries are exploiting numerical abilities, including in some cases the inability of the models to do math.
- Both models seem to be biased towards predicting the contradiction label, with the number of predicted contradiction labels being almost 20% more than the number of entailment labels.

Attention heads: We look at the role of attention heads for both the standard evaluation setting (using MNLI-m as the target dataset) and for the adversarial evaluation setting (using ANLI-A1 as the target dataset). In (Michel et al., 2019) the authors note that at test time, many attention heads can be removed individually without impacting model performance. In Section 6.1 we noted that a BERT_{BASE} trained on MNLI data did worse than Random on the ANLI-A1 test set. Interestingly, we saw that we had to prune almost 90% of the attention heads of the model for it to match Random. Figure 1 accuracy of a BERT_{BASE} model as a function of the number of attention heads. The BERT_{BASE} model was trained on MNLI only data.

6.3 Training on multiple data sources

In line with our stated hypothesis, we now mix multiple datasets to construct our train data. The hope is that this will enable us to capture a diverse range of latent signals from multiple data generation processes. For efficiency purposes, we currently sample from the data ratio to gain an intuition on how well our models are doing. In this setting, we only report the results for the BERT_{BASE} and RoBERTa_{BASE} models

Unfortunately, so far our results are in line with the observations in Section 6.1. Both the transformer models are doing remarkably well on the Dev set, but poorly on the ANLI-A1

Test set.

6.4 Advice

Any suggestions on how to further improve the models would be highly welcome. Specifically in addressing the limitation of models with respect to their numerical abilities as noted above in Section 6.2

References

- Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*.
- Max Glockner, Vered Shwartz, and Yoav Goldberg. 2018. Breaking nli systems with sentences that require simple lexical inferences. *arXiv preprint arXiv:1805.02266*.
- R Thomas McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. *arXiv preprint arXiv:1902.01007*.
- Paul Michel, Omer Levy, and Graham Neubig. 2019. Are sixteen heads really better than one? In *Advances in Neural Information Processing Systems*, pages 14014–14024.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2019. Adversarial nli: A new benchmark for natural language understanding. *arXiv preprint arXiv:1910.14599*.
- Adina Williams, Nikita Nangia, and Samuel R Bowman. 2017. A broad-coverage challenge corpus for sentence understanding through inference. *arXiv preprint arXiv:1704.05426*.