# Textual Data Augmentation for Efficient Active Learning on Tiny Datasets

Husam Quteineh, Spyridon Samothrakis and Richard Sutcliffe
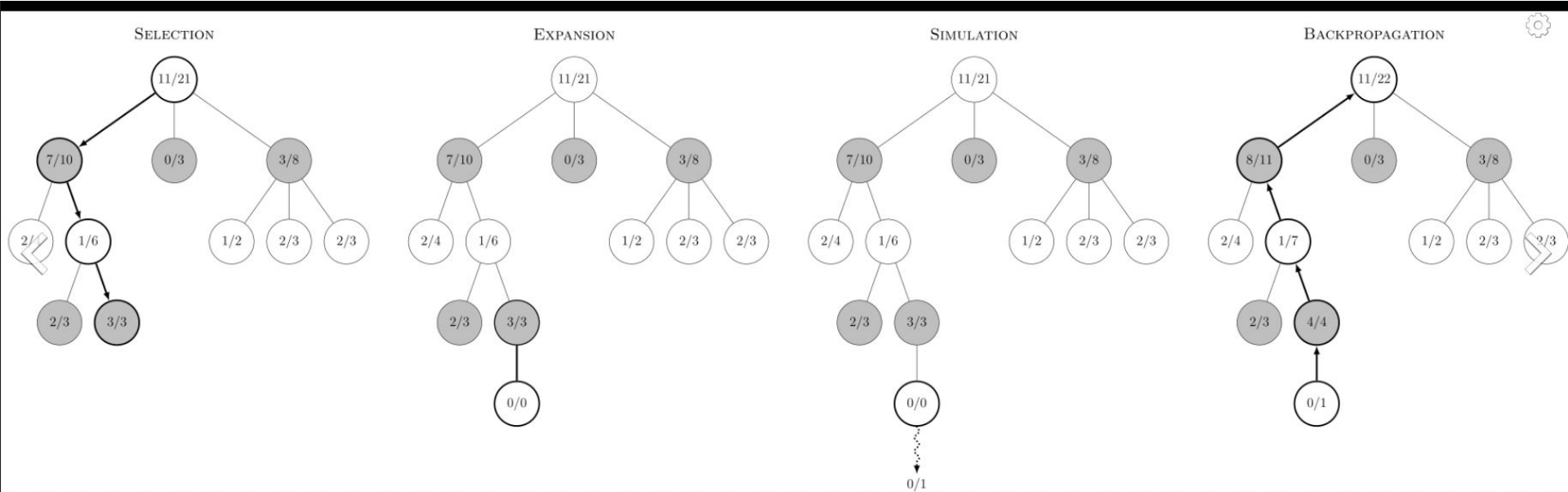
EMNLP 2020

# Problem

- Active Learning:
  - Train a classifier on a small labeled data
  - Select the most informative unlabeled data and manually label them
  - Retrain the classifier on the combination of the new labeled and the training data


- Issue: If there is not unlabeled data AL is not possible
- Solution: This paper propose to use GPT-2 to generate unlabeled data for AL
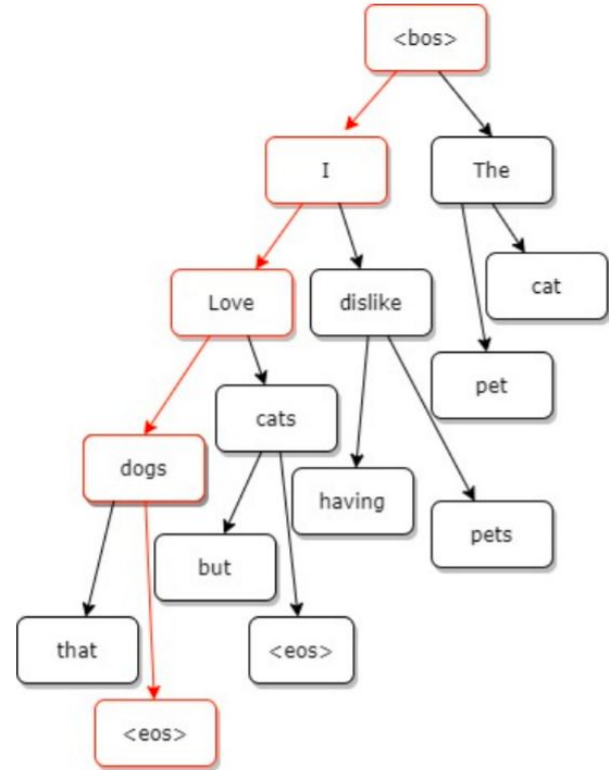  - The generation is guided by the performance of the classifier

# Solution

- Pretrain the GPT-2 model on the available labeled data
  - <BOS> W1, W2, … , Wn <EOS>
- Generate sentences and use Monte Carlo Tree Search to find the best examples
  - Model Uncertainty
  - Sentence Diversity
- Sentences with highest score are manually labeled
- New labeled data is added to training set and model is retrained

# MCTS

# Tree Search

- Each node is a word
- Children are top k nodes with highest probability based on language model
- Tree is expanded to generate multiple sentences, i.e., <EOS> nodes
- A sentence is the path from root, i.e. <BOS>, to leaf, i.e. <EOS>
- Whenever <EOS> is generate the path is evaluated

# Rewards

- Path evaluation is based on two criteria:
  - Uncertainty of the model to predict label:

$$H_n(P) = -\sum_{i=1}^{n} p_i \log_b p_i \cdot \frac{1}{\log_b n} \qquad f(x_{ent}) = \begin{cases} 0, & \text{if } x_{ent} \geq \theta_{ent} \\ x_{ent}, & \text{otherwise} \end{cases}$$

  - Difference with existing data:

$$f(x_{ent}, x_{sim}) = \begin{cases} 0, & \text{if } x_{ent} \geq \theta_{ent} \\ 0, & \text{if } x_{sim} > \theta_{sim} \\ x_{ent}, & \text{otherwise} \end{cases}$$

# Node Expansion

- A node is randomly expanded according to its score:

$$UCB = max(N_i) + C\sqrt{\frac{2 \times lnS_p}{S_i}}$$

  - Whenever a <EOS> node is generated the model evaluate the path and update every node on the path using above equation:
    - Nodes with higher reward are promoted
    - Nodes with more sentences generated from their parent are promoted

# Data Selection

- Top n sentences with highest rewards are selected from the tree
- Each selected sentence is manually labeled
- Labeled data are added to the training set and model is retrained


- Model is evaluated on two tasks:
  - Question Classification
  - Sentiment Analysis

# Results - Question Classification

| AL Run | MCTS | | NGDG |
|--------|------|--------|------|
| | Diversity | Uncert. | |
| Start | 65 (30#) | 65 (30#) | 65 (30#) |
| 1 | 68 (48#) | 78 (49#) | 78 (47#) |
| 2 | 86 (68#) | 82 (52#) | 86 (61#) |
| 3 | 92 (73#) | 87 (55#) | 87 (72#) |
| 4 | 91 (76#) | 89 (59#) | 88 (83#) |
| 5 | 92 (83#) | 91 (71#) | 86 (89#) |
| 6 | 91 (91#) | 90 (76#) | 84 (103#) |
| 7 | 90 (94#) | 89 (87#) | 84 (113#) |
| 8 | 91 (98#) | 90 (94#) | 88 (126#) |

# Results - Sentiment Analysis

| AL Run | MCTS | | NGDG |
|--------|------|------|------|
| | Diversity | Uncert. | |
| Start | 73 (20#) | 73 (20#) | 73 (20#) |
| 1 | 74 (34#) | 77 (34#) | 69 (32#) |
| 2 | 79 (41#) | 76 (44#) | 72 (43#) |
| 3 | 79 (50#) | 78 (48#) | 75 (55#) |
| 4 | 80 (60#) | 80 (54#) | 76 (79#) |
| 5 | 80 (65#) | 80 (55#) | 75 (92#) |
| 6 | 80 (79#) | 80 (62#) | 76 (103#) |
| 7 | 83 (87#) | 80 (64#) | 79 (116#) |
| 8 | 83 (95#) | 79 (69#) | 78 (124#) |

# Generated Sentences

| #  | Example |
|----|---------|
| 1  | Why did Einstein lose a fight with cancer? |
| 2  | Why did Lincole Ljungberg retire? |
| 3  | Why was Lorne L. Huntington's IQ so low? |
| 4  | What are three fundamental principles of socialism? |
| 5  | What is D.C.'s major metropolitan area? |
| 6  | When was Antarctica formed? |
| 7  | When did animals roam the earth? |
| 8  | Where can a geologist find fossils? |
| 9  | Where can an electrician find work? |
| 10 | How did Moses rule the ancient tribes? |
| 11 | How often have animals been killed by car crashes? |
| 12 | Which is Fordham's largest engineering college? |

# Thanks