

## Contents

<b>1</b>	<b>Notes from &lt;2021-02-02 Tue&gt;</b>	<b>2</b>
1.1	Present at meeting . . . . .	2
1.2	Agenda . . . . .	2
1.3	Notes . . . . .	2
<b>2</b>	<b>Goal: Publish a paper</b>	<b>2</b>
<b>3</b>	<b>Topic: NLP Data Augmentation</b>	<b>2</b>
3.1	Approach . . . . .	2
<b>4</b>	<b>Focus on Classification Tasks</b>	<b>2</b>
4.1	Ari's suggestion: . . . . .	2
4.2	Kalpesh: . . . . .	3
<b>5</b>	<b>Evaluation:</b>	<b>3</b>
<b>6</b>	<b>Dataset: IMDB</b>	<b>3</b>
<b>7</b>	<b>TODO Create a Jupyter Notebook for Baseline</b>	<b>3</b>
<b>8</b>	<b>Shared Vector Embedding space is enormous. Focus on a small subset.</b>	<b>3</b>
<b>9</b>	<b>Will discuss further after setting up toy retrieval pipeline</b>	<b>3</b>
<b>10</b>	<b>Oracle's previous approach</b>	<b>3</b>
<b>11</b>	<b>Kalpesh suggested Don't stop Pretraining embedding an additional pretraining with relevant topic</b>	<b>4</b>
<b>12</b>	<b>Get better semantic representation(embedding) by pre-training with relevant domain corpus</b>	<b>4</b>
<b>13</b>	<b>Ari suggested two papers on data augmentation strategy</b>	<b>4</b>
<b>14</b>	<b>Additional Papers</b>	<b>4</b>
14.1	VAMPIRE( <a href="https://github.com/allenai/vampire">https://github.com/allenai/vampire</a> ) . . . . .	4

<b>15 ACTIONS</b>	<b>4</b>
15.1 <b>TODO</b> Action #4 Talk to companies #4:SAM . . . .	4
15.2 <b>TODO</b> Action #5 foo #5:ME . . . . .	4

## 1 Notes from <2021-02-02 Tue>

### 1.1 Present at meeting

, - [X] Peter , - [X] Sudarshan , - [X] Kalpesh , - [X] Ari , - [X] Sweta , - [X] Mike

### 1.2 Agenda

, - Comments and corrections to last meeting notes , - Reports from the sub teams , - Discussion , - Final round

### 1.3 Notes

, . . . ,

## 2 Goal: Publish a paper

## 3 Topic: NLP Data Augmentation

### 3.1 Approach

- Find existing examples within large corpus such as Common Crawl using embedding
- Generating from scratch(May not efficient)
- Modifying existing data(Style Transfer, . . .)

## 4 Focus on Classification Tasks

### 4.1 Ari's suggestion:

- Intent Classification on banking domain (Open a bank account)
- Sentiment Classification(3 classes)

## **4.2 Kalpesh:**

- Citation intent classification
- Topic Classification

## **5 Evaluation:**

- Held-out Test Accuracy
- Adversarial held-out accuracy
- Checklist Evaluation for Linguistic Generalization
- Kalpesh: 10-15 Test Datasets for generalization - Which paper?

## **6 Dataset: IMDB**

## **7 TODO Create a Jupyter Notebook for Baseline**

- Sentiment Classification on IMDB

## **8 Shared Vector Embedding space is enormous. Focus on a small subset.**

- Linguistic phenomena
- TF-IDF
- Prefix String
- Kalpesh: Dense Embedding(DPR)

## **9 Will discuss further after setting up toy retrieval pipeline**

## **10 Oracle's previous approach**

- Collecting the farthest embedding.
- May collect irrelevant data. Thus, think about other way

- e.g. Collect task-relevant examples and take the farthest embedding

## 11 Kalpesh suggested Don't stop Pretraining embedding an additional pretraining with relevant topic

- 

## 12 Get better semantic representation(embedding) by pre-training with relevant domain corpus

- ROBERTA-DAPT and ROBERTA-TAPT are available online.

<https://huggingface.co/allenai>

## 13 Ari suggested two papers on data augmentation strategy

## 14 Additional Papers

### 14.1 VAMPIRE(<https://github.com/allenai/vampire>)

- This is the augmentation method for the previous paper "don't stop pretraining"
- Got it working without any issue.

## 15 ACTIONS

, This is the general list of Actions

15.1 TODO Action #4 Talk to companies #4:SAM

15.2 TODO Action #5 foo #5:ME