

adversarial training methods for semi-supervised text classification

abstract

adversarial training + virtual
adversarial training



noise in neural network
input 자체에는 적용X

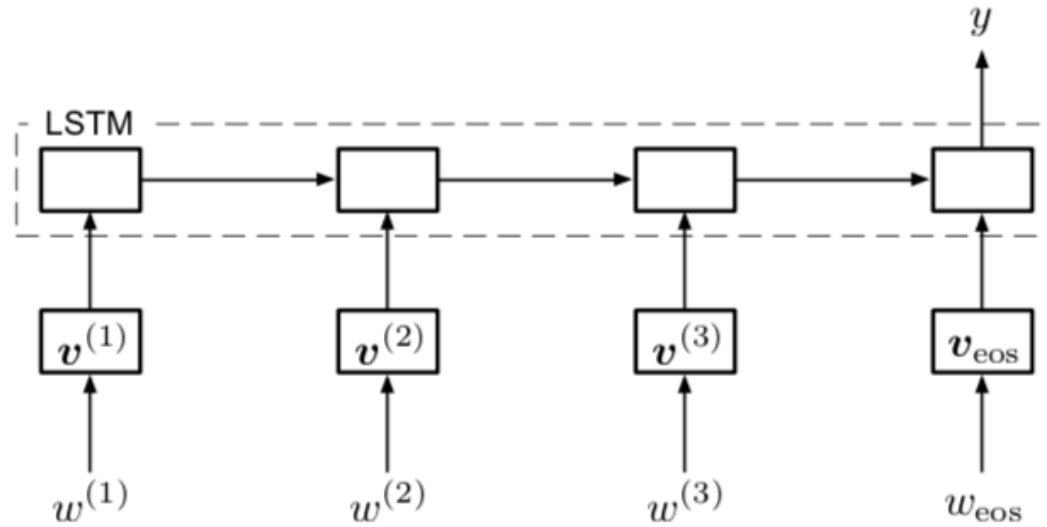
introduction

- adversarial training
 - origin sample / noise sample을 모두 정확하게 구분하는 model 만들기 위한 과정
 - label 필수
- virtual adversarial training
 - unlabeled sample
 - 모델의 regularization
 - > origin sample & noise sample 모두 같은 출력

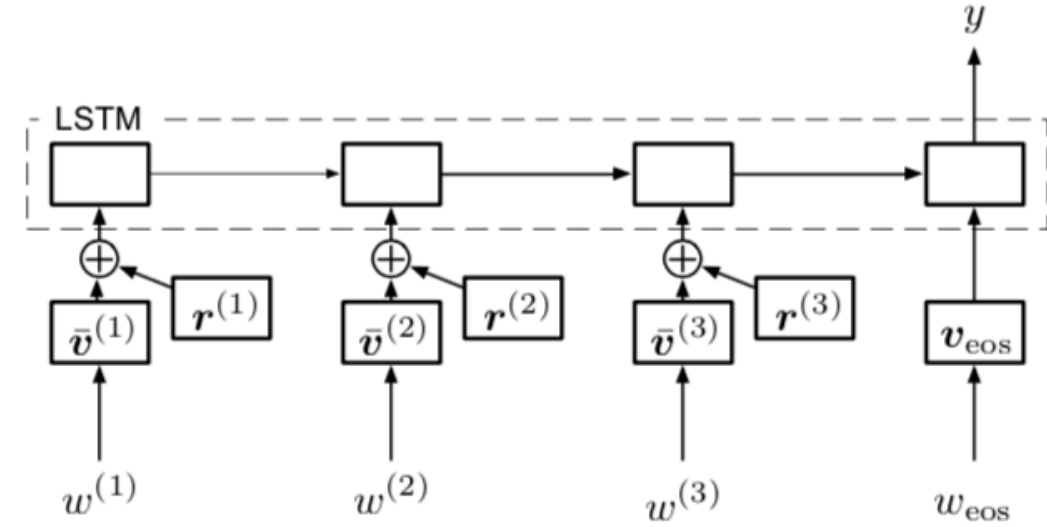
악의적 input에 대한 방어로 작용

word embedding에 접근할 수 없기 때문에
classifier regularization으로서 분류 기능 안정화를 제시

method



(a) LSTM-based text classification model.



(b) The model with perturbed embeddings.

$$\bar{v}_k = \frac{\mathbf{v}_k - \mathbb{E}(\mathbf{v})}{\sqrt{\text{Var}(\mathbf{v})}} \text{ where } \mathbb{E}(\mathbf{v}) = \sum_{j=1}^K f_j \mathbf{v}_j, \text{Var}(\mathbf{v}) = \sum_{j=1}^K f_j (\mathbf{v}_j - \mathbb{E}(\mathbf{v}))^2$$

method

adversarial training / loss

$$-\log p(y \mid \mathbf{x} + \mathbf{r}_{\text{adv}}; \boldsymbol{\theta}) \text{ where } \mathbf{r}_{\text{adv}} = \arg \min_{\mathbf{r}, \|\mathbf{r}\| \leq \epsilon} \log p(y \mid \mathbf{x} + \mathbf{r}; \hat{\boldsymbol{\theta}})$$

$$\mathbf{r}_{\text{adv}} = -\epsilon \mathbf{g} / \|\mathbf{g}\|_2 \text{ where } \mathbf{g} = \nabla_{\mathbf{x}} \log p(y \mid \mathbf{x}; \hat{\boldsymbol{\theta}}).$$

$$L_{\text{adv}}(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{n=1}^N \log p(y_n \mid \mathbf{s}_n + \mathbf{r}_{\text{adv},n}; \boldsymbol{\theta})$$

virtual adversarial training / loss

$$\text{KL}[p(\cdot \mid \mathbf{x}; \hat{\boldsymbol{\theta}}) \parallel p(\cdot \mid \mathbf{x} + \mathbf{r}_{\text{v-adv}}; \boldsymbol{\theta})]$$

$$\text{where } \mathbf{r}_{\text{v-adv}} = \arg \max_{\mathbf{r}, \|\mathbf{r}\| \leq \epsilon} \text{KL}[p(\cdot \mid \mathbf{x}; \hat{\boldsymbol{\theta}}) \parallel p(\cdot \mid \mathbf{x} + \mathbf{r}; \boldsymbol{\theta})]$$

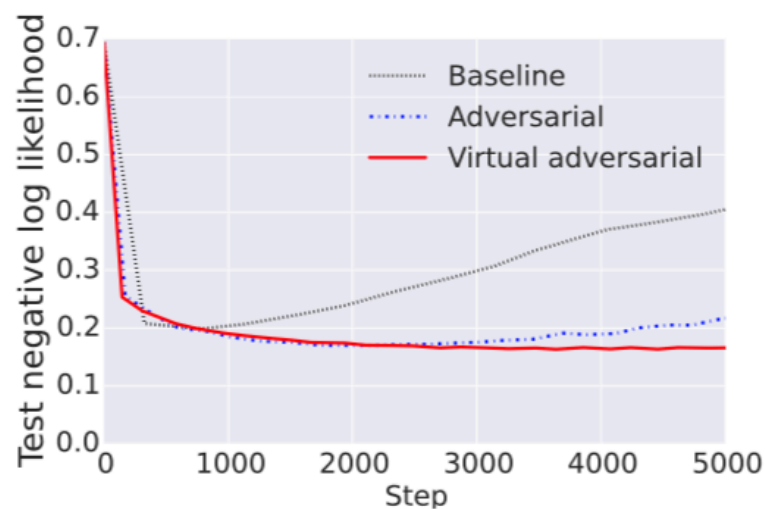
$$\mathbf{r}_{\text{v-adv}} = \epsilon \mathbf{g} / \|\mathbf{g}\|_2 \text{ where } \mathbf{g} = \nabla_{\mathbf{s} + \mathbf{d}} \text{KL} \left[p(\cdot \mid \mathbf{s}; \hat{\boldsymbol{\theta}}) \parallel p(\cdot \mid \mathbf{s} + \mathbf{d}; \boldsymbol{\theta}) \right]$$

$$L_{\text{v-adv}}(\boldsymbol{\theta}) = \frac{1}{N'} \sum_{n'=1}^{N'} \text{KL} \left[p(\cdot \mid \mathbf{s}_{n'}; \hat{\boldsymbol{\theta}}) \parallel p(\cdot \mid \mathbf{s}_{n'} + \mathbf{r}_{\text{v-adv},n'}; \boldsymbol{\theta}) \right]$$

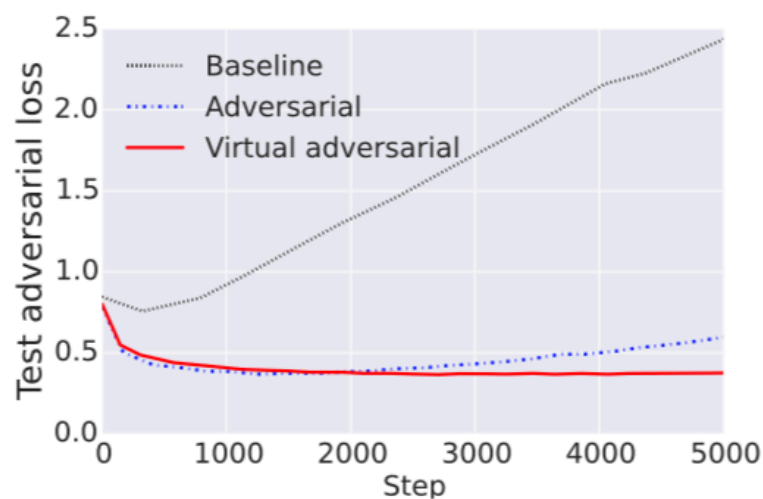
experiment

Table 1: Summary of datasets. Note that unlabeled examples for the Rotten Tomatoes dataset are not provided so we instead use the unlabeled Amazon reviews dataset.

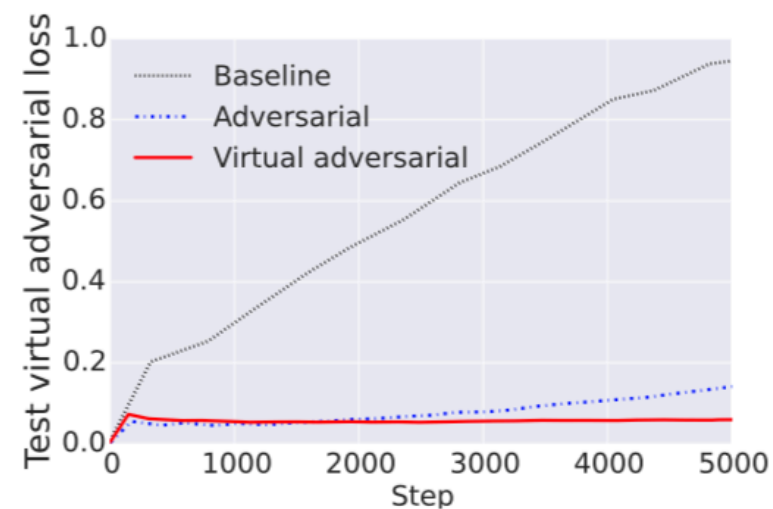
| | Classes | Train | Test | Unlabeled | Avg. T | Max T |
|-----------------|---------|---------|--------|-----------|----------|---------|
| IMDB | 2 | 25,000 | 25,000 | 50,000 | 239 | 2,506 |
| Elec | 2 | 24,792 | 24,897 | 197,025 | 110 | 5,123 |
| Rotten Tomatoes | 2 | 9596 | 1066 | 7,911,684 | 20 | 54 |
| DBpedia | 14 | 560,000 | 70,000 | – | 49 | 953 |
| RCV1 | 55 | 15,564 | 49,838 | 668,640 | 153 | 9,852 |



(a) Negative log likelihood



(b) $L_{\text{adv}}(\theta)$



(c) $L_{\text{v-adv}}(\theta)$

conclusion

- classification, word embedding에 뛰어난 성과
- 음성, 비디오와 같은 순차적 작업에 적용 가능성

code : https://github.com/tensorflow/models/tree/master/adversarial_text