

Adversarial Training with Contrastive Learning in NLP

Procesamiento de Lenguaje Natural

Francisco Javier Sáez Maldonado

14 de marzo de 2022

Máster en Ciencia de Datos

*Escuela Politécnica Superior
Universidad Autónoma de Madrid*

Herramientas

Adversarial Training

Contrastive Learning

Framework

Experimentos y resultados

- **Tarea:** Modelado del lenguaje (LM) y Traducción automática (NMT)

- **Tarea:** Modelado del lenguaje (LM) y Traducción automática (NMT)
- **Objetivo:** Conseguir modelos que sean más robustos semánticamente:

Inputs parecidos \implies Outputs parecidos

Herramientas

Herramientas

Adversarial Training

Definición (Adversarial Learning)

Técnica usada en el aprendizaje automático para, usando información sobre un modelo, crear ataques maliciosos para causar fallos en el modelo

Definición (Adversarial Learning)

Técnica usada en el aprendizaje automático para, usando información sobre un modelo, crear ataques maliciosos para causar fallos en el modelo

Definición (Adversarial example)

Ejemplo diseñado para engañar al modelo, creado introduciendo una perturbación en un ejemplo original.

¿ Cómo ayuda el aprendizaje adversario a nuestros modelos ?

Técnicas de creación ejemplos adversarios en NLP

- Visuales: (Morris u. a., 2020).

Original Input	This film has a special place in my heart	Positive
Adversarial example	This film has a special plcae in my herat	Negative

Cuadro 1: Ejemplo extraído de (Gao u. a., 2018)

- Semánticas:(Jin u. a., 2019)

Original
<u>Perfect</u> performance by the actor → Positive (99%)
.....
Adversarial
<u>Spotless</u> performance by the actor → Negative (100%)

Figura 1: Ejemplo extraído de (Jin u. a., 2019)

Ejemplos adversarios

Dada una secuencia $s = \{x_1, \dots, x_T\}$ de tokens

1. Creamos una representación embebida en un espacio continuo

$$\mathbf{E}x_i = e_i.$$

2. Añadimos una pequeña perturbación en el embedding

$$e'_i = e_i - \epsilon \frac{g}{\|g\|_2},$$

siendo $g = \nabla_{e_i} J(s, \theta)$ y J la función de coste.

Función de coste actual:

$$\mathcal{J}(\theta) = \sum_s \mathcal{L}(s, \theta) + \alpha \sum_{s'} \mathcal{L}_{adv}(s', \theta), \quad \alpha \in [0, 1].$$

Herramientas

Contrastive Learning

Contrastive Learning

Idea: Acercar las representaciones de ejemplos positivos (de la misma clase) y alejar las de los ejemplos negativos (resto de ejemplos).

Ejemplo

Original: Elefante. Positivo: Hipopótamo. Negativo: Pistola.

Definición (Pérdida contrastiva)

Sean a_i las entradas originales, p_{a_i} ejemplos positivos y n_a ejemplos negativos. Se define la pérdida contrastiva como:

$$\mathcal{L}_{cont} = - \sum_{a_i \in A} \log \frac{\exp(a_i \cdot p_{a_i} / \tau)}{\sum_{n_a \in A - \{a_i\}} \exp(a_i \cdot n_a / \tau)}$$

Framework

Consideramos :

- $\mathbf{S} = \{s_1, \dots, s_B\}$ un conjunto de frases, donde cada frase $s_k = \{x_{k1}, \dots, x_{kN}\}$ tiene N tokens.

Consideramos :

- $\mathbf{S} = \{s_1, \dots, s_B\}$ un conjunto de frases, donde cada frase $s_k = \{x_{k1}, \dots, x_{kN}\}$ tiene N tokens.
- $\mathbf{E}s_k = \{\mathbf{E}x_{k1}, \dots, \mathbf{E}x_{kN}\} = \{e_{k1}, \dots, e_{kN}\}$ el embedding al espacio continuo de cada frase.

Consideramos :

- $\mathbf{S} = \{s_1, \dots, s_B\}$ un conjunto de frases, donde cada frase $s_k = \{x_{k1}, \dots, x_{kN}\}$ tiene N tokens.
- $\mathbf{E}s_k = \{\mathbf{E}x_{k1}, \dots, \mathbf{E}x_{kN}\} = \{e_{k1}, \dots, e_{kN}\}$ el embedding al espacio continuo de cada frase.
- El vocabulario \mathcal{V} y un subconjunto **restringido** del mismo del que excluimos palabras incompletas (caracteres individuales o símbolos) \mathcal{V}_R

Consideramos :

- $\mathbf{S} = \{s_1, \dots, s_B\}$ un conjunto de frases, donde cada frase $s_k = \{x_{k1}, \dots, x_{kN}\}$ tiene N tokens.
- $\mathbf{E}s_k = \{\mathbf{E}x_{k1}, \dots, \mathbf{E}x_{kN}\} = \{e_{k1}, \dots, e_{kN}\}$ el embedding al espacio continuo de cada frase.
- El vocabulario \mathcal{V} y un subconjunto **restringido** del mismo del que excluimos palabras incompletas (caracteres individuales o símbolos) \mathcal{V}_R
- La función que nos da la restricción a \mathcal{V}_R :

$$\mathcal{M}(\mathbf{E}x_{ki}) = \begin{cases} 1 & \text{if } x_{ki} \in \mathcal{V}_R \\ 0 & \text{otherwise} \end{cases}$$

(evita tomar candidatos a modificación adversaria sin sentido).

Consideramos :

- $\mathbf{S} = \{s_1, \dots, s_B\}$ un conjunto de frases, donde cada frase $s_k = \{x_{k1}, \dots, x_{kN}\}$ tiene N tokens.
- $\mathbf{E}s_k = \{\mathbf{E}x_{k1}, \dots, \mathbf{E}x_{kN}\} = \{e_{k1}, \dots, e_{kN}\}$ el embedding al espacio continuo de cada frase.
- El vocabulario \mathcal{V} y un subconjunto **restringido** del mismo del que excluimos palabras incompletas (caracteres individuales o símbolos) \mathcal{V}_R
- La función que nos da la restricción a \mathcal{V}_R :

$$\mathcal{M}(\mathbf{E}x_{ki}) = \begin{cases} 1 & \text{if } x_{ki} \in \mathcal{V}_R \\ 0 & \text{otherwise} \end{cases}$$

(evita tomar candidatos a modificación adversaria sin sentido).

- Llamaremos h_{kj} a la representación obtenida de la frase s_k .

Adversarial Training with Contrastive Learning (ATCL)

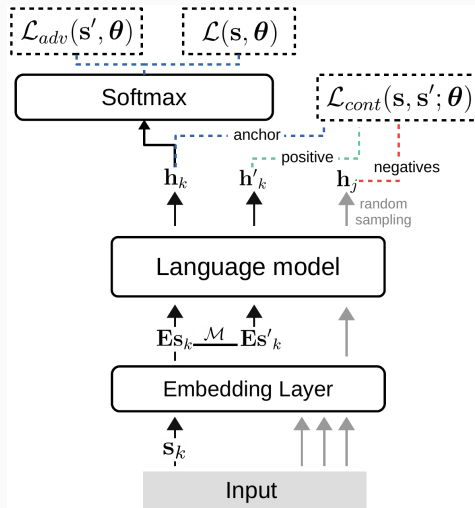


Figure 2: ACTL Framework. Image obtained from the original paper (Rim u. a., 2021).

Consideraciones sobre ATCL

- Con el ejemplo adversario alejamos la representación de la frase original, y forzamos la cercanía mediante contrastive learning.
- Para la función de pérdida contrastiva, $\mathcal{L}_{\text{cont}}$, los ejemplos negativos se samplean del conjunto $\mathbf{H}(\mathbf{S}) - \{h_{kj}\}$, para evitar tomar a la misma h_{kj} como ejemplo negativo.

Definición (Función de pérdida de ATCL)

La función de pérdida que se utiliza para entrenar el framework ATCL es:

$$\mathcal{J}(\theta)_{ATCL} = \sum_{\mathbf{s}, \mathbf{s}'} (\mathcal{L} + \alpha \mathcal{L}_{adv} + \beta \mathcal{L}_{cont})$$

Experimentos y resultados

Datasets:

- omega Yea

Gracias por su atención

Referencias

- [Gao u. a. 2018] GAO, Ji ; LANCHANTIN, Jack ; SOFFA, Mary L. ; QI, Yanjun: Black-box Generation of Adversarial Text Sequences to Evade Deep Learning Classifiers. (2018). – URL <http://arxiv.org/abs/1801.04354>
- [Jin u. a. 2019] JIN, Di ; JIN, Zhijing ; ZHOU, Joey T. ; SZOLOVITS, Peter: Is BERT Really Robust? Natural Language Attack on Text Classification and Entailment. (2019). – URL <http://arxiv.org/abs/1907.11932>
- [Morris u. a. 2020] MORRIS, John X. ; LIFLAND, Eli ; YOO, Jin Y. ; QI, Yanjun: TextAttack: A Framework for Adversarial Attacks in Natural Language Processing. (2020). – URL <https://arxiv.org/abs/2005.05909>
- [Rim u. a. 2021] RIM, Daniela N. ; HEO, DongNyeong ; CHOI, Heeyoul: Adversarial Training with Contrastive Learning in NLP. (2021). – URL <https://arxiv.org/abs/2109.09075>