# GAN-BERT: Generative Adversarial Learning for Robust Text Classification

Lecture: Seminar Klassifikation und Clustering Lecturer: PD Dr. Stefan Langer Presenters: Mengmeng Xu, Bolei Ma, Zheyu Zhang, Han Yang

Date: 23.01.2023



# Agenda

- 1. Introduction
- 2. Basic Models
  - 2.1 BERT
  - 2.2 GANs
- 3. SS-GANs
- 4. GAN-BERT
- 5. Experimental Settings
- 6. Results & Analysis
- 7. Conclusion



# Introduction

- Goal:
  - 1. Save time on data annotation (reduce the scale of labeled examples)
  - 2. Maintain good performances
- Main Idea:

Use GAN to extend the fine-tuning of BERT-like architectures with unlabeled data in a generative adversarial setting

### Main Contributions:

Extended the limits of Transformer-based architectures (i.e., BERT) in poor training conditions
 Systematically improved the robustness of such architectures, while not introducing additional costs to the inference



# **BERT** – Bidirectional Encoder Representations from Transformers

- A stack of Transformer Encoder Layers that consists of multiple self-attention "heads"
- Workflow (crucial part):
  - 1. Pre-training: MLM (masked language modeling) & NSP (next sentence prediction) Special tokens in NSP: [CLS], [SEP]
  - 2. Fine-tuning: (basically tells BERT what to ignore)
    - Presumably teach the model to rely more on the representations useful for the task at hand
    - Possibilities of improving the fine-tuning of BERT:
      - 1. Taking more layers into account: learning a complementary representation of the information in deep and output layers, using a weighted combination of all layers instead of the final one and layer dropout
      - 2. **Two-stage fine tuning:** introduces an intermediate supervised training stage between pre-training and fine-tuning. Ben-David et al. (2020) propose a pivot-based variant of MLM to fine-tune BERT for domain adaptation
      - 3. Adversarial token perturbations: improve the robustness of the model (Zhu et al., 2019)
      - 4. Adversarial regularization: helps alleviate pre-trained knowledge forgetting and therefore prevents BERT from overfitting to downstream tasks (Jiang et al., 2019a)
      - 5. **Mixout regularization:** improves the stability of BERT fine-tuning even for a small number of training examples (Lee et al., 2019)



### **BERT** – Bidirectional Encoder Representations from Transformers

Using different open-source models on the Hugging Face for fine-tuning

∋ BERT	🔞 SIGN IN 🛛 🖋 MODELS 🗢 FORU
Overview	Docs »BERT View page source
BertConfig	
BertTokenizer	
BertTokenizerFast	REDT
Bert specific outputs	DLINI
BertModel	
BertForPreTraining	
BertModelLMHeadModel	Overview
BertForMaskedLM	
BertForNextSentencePrediction	The BERT model was proposed in BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding by Jacob Devlin, Ming-Wei Chang, Kenton Lee and Kristina Toutanova. It's a bidirectional transformer pre-trained using a combination of masked language modeling objective and next
BertForSequenceClassification	sentence prediction on a large corpus comprising the Toronto Book Corpus and Wikipedia.
BertForMultipleChoice	
BertForTokenClassification	The abstract from the paper is the following:
BertForQuestionAnswering	We introduce a new language representation model called BERT, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent
TFBertModel	language representation models, BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and
TFBertForPreTraining	right context in all layers. As a result, the pre-trained BER I model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide rance of tasks, such as question answering and language inference, without substantial task-specific architecture modifications.
TFBertModelLMHeadModel	
TFBertForMaskedLM	BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven natural language processing tasks, including pushing
TFBertForNextSentencePrediction	the GLUE score to 60.5% (1.7% point absolute improvement), multinul accuracy to 86.7% (4.6% absolute improvement), SQUAD v1.1 question answering Test F1 to 93.2 (1.5 point absolute improvement) and SQUAD v2.0 Test F1 to 83.1 (5.1 point absolute improvement).
TFBertForSequenceClassification	
TFBertForMultipleChoice	Tips:

#### **Bert Model Overview on Hugging Face**

#### **Bert Model Variations**

Model	#params	Language
bert-base-uncased	110M	English
<u>bert-large-uncased</u>	340M	English
<u>bert-base-cased</u>	110M	English
<u>bert-large-cased</u>	340M	English
<u>bert-base-chinese</u>	110M	Chinese
<u>bert-base-multilingual-cased</u>	110M	Multiple
bert-large-uncased-whole-word-masking	340M	English
bert-large-cased-whole-word-masking	340M	English

src: https://huggingface.co/bert-base-uncased

src: https://huggingface.co/transformers/v3.0.2/model\_doc/bert.html



# GANs – Generative Adversarial Networks

- An approach to generative modeling using deep learning methods, such as convolutional neural networks, applied in CV
- Automatically discovering and learning the regularities or patterns in input data
- Generate or output new examples that plausibly could have been drawn from the original dataset

### **Components of GANs:**

- > Generator (G): be trained to generate new examples
- > Discriminator (D): tries to classify examples as either real (from the domain) or fake (generated)
- Process: The two models are trained together in a zero-sum game, adversarial, until the discriminator model is fooled about half the time, meaning the generator model is generating plausible examples



### **GANs – Generative Adversarial Networks**





# **Detailed Illustration of GANs**

### **Generator Model**



- 1. Takes a fixed-length random vector as input and generates a sample in the domain
- 2. The vector is drawn randomly from a Gaussian distribution, and the vector is used to seed the generative process
- 3. Points in this multidimensional vector space will correspond to points in the problem domain after training, forming a compressed representation of the data distribution



# **Detailed Illustration of GANs**

### **Discriminator Model**



- 1. The discriminator model takes an example from the domain as input (real or generated) and predicts a binary class label of real or fake (generated)
- 2. The discriminator is a normal (and well understood) classification model

### **Adversarial Process**



Example of the Generative Adversarial Network Model Architecture

- SS-GANs: semi-supervised learning in a GAN framework
- A (k + 1)-class objective:
- "real" examples are classified into target classes [1, 2, ... , k]
- > generated / "fake" examples are classified into k + 1



• Components:

Generator (G)

> Discriminator (D)

000

 $\succ$  Loss of SS-GANs:  $L_D$  and  $L_G$ 

> We will have two Loss functions, one for L(D) and one for L(G):

$$\succ L_D = L_{D_{sup.}} + L_{D_{unsup.}}$$

$$> L_G = L_{G_{unsup.}} + L_{G_{feature matching}}$$



• First Loss of SS-GANs:  $L_D = L_{D_{sup.}} + L_{D_{unsup.}}$ 

 $> L_{D_{sup.}}$  measures the error in assigning the wrong class to a real example among the original k categories.

 $> L_{D_{unsup.}}$  measures the error in incorrectly recognizing a real (unlabeled) example as fake and not recognizing a fake example.

$$\begin{split} L_{\mathcal{D}_{\text{sup.}}} &= -\mathbb{E}_{x, y \sim p_d} \log \left[ p_{\text{m}}(\hat{y} = y | x, y \in (1, ..., k)) \right] \\ & \text{real data} \\ L_{\mathcal{D}_{\text{unsup.}}} &= -\mathbb{E}_{x \sim p_d} \log \left[ 1 - p_{\text{m}} \left( \hat{y} = y | x, y = k + 1 \right) \right] \\ & \text{real data} \\ & - \mathbb{E}_{x \sim \mathcal{G}} \log \left[ p_{\text{m}}(\hat{y} = y | x, y = k + 1) \right] \\ & \text{fake data} \end{split}$$

Probs of associating real data in k classes

Probs of associating fake data into k+1 class

### **SS-GANs**

• Second Loss of SS-GANs:  $L_G = L_{G_{unsup.}} + L_{G_{feature matching}}$ 

> The  $L_{G_{unsup.}}$  considers the error induced by fake examples correctly identified by D

$$L_{\mathcal{G}_{unsup.}} = -\mathbb{E}_{x \sim \mathcal{G}} \log \left[1 - p_m \left(\hat{y} = y | x, y = k+1\right)\right]$$
fake data

Probs of associating real data in k classes

Probs of associating fake data into k+1 classes

• Second Loss of SS-GANs:  $L_G = L_{G_{unsup}} + L_{G_{feature matching}}$ 

> Reason for feature matching: G is expected to generate examples that are similar to the ones sampled from the real distribution P(d).

> Feature matching loss calculation: Let f(x) denote the activation on an intermediate layer of D; the generator should produce examples whose intermediate representations provided in input to D are very similar to the real ones.

$$L_{G_{\text{feature matching}} = \left\| \mathbb{E}_{x \ \sim \ p_{d}} f(x) - \mathbb{E}_{x \ \sim \ \mathcal{G}} f(x) \right\|_{2}^{2}$$

real data distribution in input to D

# **GAN-BERT**



- Two components:
- > task-specific layers, as in the usual BERT fine-tuning
- > SS-GAN layers to enable semi-supervised learning



0 0 C

# **GAN-BERT**

### • GAN-BERT architecture:

> G generates a set of fake examples F given a random distribution.

Fake examples F, along with unlabeled U and labeled L vector representations computed by BERT are used as input for the discriminator D.



### **GAN-BERT**

- Loss of GAN-BERT (as of SS-GANs):
- $\succ L_D = L_{D_{sup.}} + L_{D_{unsup.}}$
- $\succ L_G = L_{G_{unsup.}} + L_{G_{feature matching}}$
- Weight update: G, D, BERT



000

# **GAN-BERT**

- Pre-trained BERT model:



Given input sentence:

 $s=\left( t_{1},...,t_{n}
ight)$ 

- **BERT** output n+2 vector representations:

 $\left(h_{CLS},h_{t_1},...,h_{t_n},h_{SEP}
ight)$ 

- We adopt the *h<sub>cls</sub>* representation as a sentence embedding









- MLP with one hidden layer
- Activated by a leaky-relu function









logits = self.logit(last\_rep)
probs = self.softmax(logits)

return last\_rep, logits, probs

- Followed by a softmax layer for the final prediction

23/01/2023

# **GAN-BERT**

### **Training Procedure:**

23/01/2023

# Encode real data in the Transformer model\_outputs = transformer(b\_input\_ids, attention\_mask=b\_input\_mask) 1. Encode real data in Transformer hidden\_states = model\_outputs[-1] # Generate fake data that should have the same distribution of the ones # encoded by the transformer. # First noisy input are used in input to the Generator noise = torch.zeros(real\_batch\_size, noise\_size, device=device).uniform\_(0, 1) 2. Generate fake data # Gnerate Fake data gen\_rep = generator(noise) # Generate the output of the Discriminator for real and fake data. # First, we put together the output of the tranformer and the generator disciminator\_input = torch.cat([hidden\_states, gen\_rep], dim=0) 3. Put real and fake data into Discriminator # Then, we select the output of the disciminator features, logits, probs = discriminator(disciminator\_input) # Finally, we separate the discriminator's output for the real and fake # data features\_list = torch.split(features, real\_batch\_size) D\_real\_features = features\_list[0] D\_fake\_features = features\_list[1] logits\_list = torch.split(logits, real\_batch\_size) D\_real\_logits = logits\_list[0] 4. Separate the Discriminator's outpu for real and fake data D\_fake\_logits = logits\_list[1] probs\_list = torch.split(probs, real\_batch\_size) D\_real\_probs = probs\_list[0] D fake probs = probs list[1]

 $cal data \\ cal data$ 

21

# **GAN-BERT**

#### Loss of GAN-BERT (as of SS-GANs):

$$L_G = L_{G_{unsup.}} + L_{G_{feature matching}}$$
 -

$$L_D = L_{D_{sup.}} + L_{D_{unsup.}} -$$

#### Weight update: G, D and BERT



# Generator's LOSS estimation g\_loss\_d = -1 \* torch.mean(torch.log(1 - D\_fake\_probs[:, -1] + self.epsilon)) g\_feat\_reg = torch.mean( torch.pow(torch.mean(D\_real\_features, dim=0) - torch.mean(D\_fake\_features, dim=0), 2)) g\_loss = g\_loss\_d + g\_feat\_reg # Disciminator's LOSS estimation logits = D\_real\_logits[:, 0:-1] log\_probs = F.log\_softmax(logits, dim=-1) # The discriminator provides an output for labeled and unlabeled real data # so the loss evaluated for unlabeled data is ignored (masked) label2one\_hot = torch.nn.functional.one\_hot(b\_labels, len(label\_list)) per\_example\_loss = -torch.sum(label2one\_hot \* log\_probs, dim=-1) per\_example\_loss = torch.masked\_select(per\_example\_loss, b\_label\_mask.to(self.device)) labeled\_example\_count = per\_example\_loss.type(torch.float32).numel() # It may be the case that a batch does not contain labeled examples, *#* so the "supervised loss" in this case is not evaluated if labeled\_example\_count == 0: D\_L\_Supervised = 0 else:

D\_L\_Supervised = torch.div(torch.sum(per\_example\_loss.to(self.device)), labeled\_example\_count)

D\_L\_unsupervised1U = -1 \* torch.mean(torch.log(1 - D\_real\_probs[:, -1] + self.epsilon))
D\_L\_unsupervised2U = -1 \* torch.mean(torch.log(D\_fake\_probs[:, -1] + self.epsilon))
d\_loss = D\_L\_Supervised + D\_L\_unsupervised1U + D\_L\_unsupervised2U

# **Experimental Part**

Goal:

- Assess the impact of GAN-BERT over sentence classification tasks characterized by different training conditions, i.e., number of examples and number of categories.
  - Binary Classification (Sentiment Analysis)
  - Multi-Class Classification (Author Identification)

Dataset:

- Sentiment (40,000 texts for Train set / 10,000 texts for Test set)
- Letters (*39,077* texts for Train set / *4,881* texts for Test set)

**Baseline**:

- **BERT-base** model fine-tuned as described in [Devlin et al., 2019] on the available training material Hyperparameters:

23

Repeat the training of each model with an increasing set of annotated material (*L*)







Dataset: Sentiment Label: Sentiment (Negative / Positive) 40, 000 training data 10, 000 testing data

- Training set:

	Text	Sentiment
count	40000	40000
unique	39745	2

- Neg: 20015
- Pos: 19985

- Test set:

	Text	Sentiment
count	10000	10000
unique	9983	2

- Neg: 4985

- Pos: 5015



### Annotated Data: 0.1% (40 texts) BERT: BERT-base Epochs: 10

-	precision	recall	f1-score	support
negative	0.640	0.580	0.608	4985
positive	0.618	0.676	0.646	5015
accuracy			0.628	10000
macro avg	0.629	0.628	0.627	10000
weighted avg	0.629	0.628	0.627	10000





### Annotated Data: 0.25% (100 texts) BERT: BERT-base Epochs: 10

-	precision	recall	f1-score	support
negative	0.704	0.647	0.674	4985
positive	0.675	0.730	0.701	5015
accuracy			0.689	10000
macro avg	0.690	0.688	0.688	10000
weighted avg	0.690	0.689	0.688	10000





### Annotated Data: 0.5% (200 texts) BERT: BERT-base Epochs: 10

-	precision	recall	fl-score	support
negative positive	0.678 0.711	0.732 0.655	0.704 0.682	4985 5015
accuracy			0.693	10000
macro avg	0.695	0.693	0.693	10000
weighted avg	0.695	0.693	0.693	10000





### Annotated Data: 1% (400 texts) BERT: BERT-base Epochs: 10

-	precision	recall	fl-score	support
negative	0.701	0.780	0.739	4985
positive	0.754	0.670	0.709	5015
accuracy			0.725	10000
macro avg	0.728	0.725	0.724	10000
weighted avg	0.728	0.725	0.724	10000



### Results

# **GAN-BERT** performs better than **BERT**:

with 0.1% and 0.25% of labeled training data



Experiment 1: Sentiment Analysis





Dataset: Letters Label: Authors **39, 077** training data **4, 881** testing data

#### - Training set:









f1-score

0.000

0.000

0.936

0.430

0.000

0.816

0.604

0.676

0.398

0.628

recall

0.000

0.000

0.949

0.422

0.000

0.806

1.000

0.454

0.676

### Annotated Data: 0.1% (39 texts) BERT: mBERT Epochs: 3

Franz Kafka

Henrik Ibsen

Virginia Woolf

Wilhelm Busch

weighted avg

accuracy macro avg

James Joyce

Friedrich Schiller

Johann Wolfgang von Goethe

precision

0.000

0.000

0.924

0.438

0.000

0.826

0.433

0.374

0.608

#### Normalized Confusion Matrix (0.1%)





### Normalized Confusion Matrix (0.25%)







23/01/2023



### Annotated Data: 1% (390 texts) BERT: mBERT Epochs: 3

Franz Kafka

Henrik Ibsen

Virginia Woolf

Wilhelm Busch

weighted avg

accuracy

macro avg

James Joyce

Friedrich Schiller

Johann Wolfgang von Goethe

### Normalized Confusion Matrix (1%)







### Normalized Confusion Matrix (5%)

23/01/2023

### Results

# **GAN-BERT** performs better than **BERT**:

with 0.1% (39 cases), 0.25% (97 cases), 0.5% (195 cases), 1% (390 cases), 5% (1953 cases), 20% of Dataset

### Experiment 2: Author Identification (Letters)





### Conclusion

- GAN-BERT improved the robustness of fine-tuning
  - even in poor training condition (low training data)
  - (intuitively) GAN provide massive negative samples during training
- GAN-BERT does not introduce additional cost during inference
  - Only additional discriminator (simple)



• More improvement of GAN-BERT is more significant on **Multiclass** Classification than **Binary** Classification

# Conclusion



#### result of GAN-BERT on 5% of dataset

# Conclusion

- BERT
- GAN (Generative Adversarial Networks)
- GAN-BERT: Fine-Tune of BERT with GAN
- Experiments
  - Sentiment Analysis
  - Author Identification
- Advantage of GAN-BERT
  - GAN provides a large number of negative samples
  - Robustness with low training data
  - Not introducing additional costs during inference
- Limitation
  - Carefully separate dataset



### Reference

- Danilo Croce, Giuseppe Castellucci, and Roberto Basili. 2020. GAN-BERT: Generative Adversarial Learning for Robust Text Classification with a Bunch of Labeled Examples. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 2114–2119, Online. Association for Computational Linguistics.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, Advances in Neural Information Processing Systems 27, pages 2672–2680. Curran Associates, Inc.
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and Xi Chen. 2016. Improved techniques for training gans. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, Advances in Neural Information Processing Systems 29, pages 2234–2242. Curran Associates, Inc.
- Anna Rogers, Olga Kovaleva and Anna Rumshisky. A Primer in BERTology: What We Know About How BERT Works. Transactions of the Association for Computational Linguistics (2020) 8: 842–866. January 01.2021.
- Jason Browniee. A Gentle Introduction to Generative Adversarial Networks(GANs).<u>Generative Adversarial Networks</u>. June 17, 2019.