

# 학습된 NarrativeKoGPT2을 이용한 Text Generation

## 1. Google Drive 연동

- 모델 파일과 학습 데이터가 저장 되어있는 구글 드라이브의 디렉토리와 Colab을 연동.

### 1.1 Google Drive 연동



아래 코드를 실행후 나오는 URL을 클릭하여 나오는 인증 코드 입력

```
In [0]: from google.colab import drive
drive.mount('/content/drive')
```

Go to this URL in a browser: [https://accounts.google.com/o/oauth2/auth?client\\_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect\\_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response\\_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly](https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly)

Enter your authorization code:  
.....

Mounted at /content/drive

**Colab 디렉토리 아래 NarrativeKoGPT2 경로 확인**

```
In [0]: !ls drive/'My Drive'/'Colab Notebooks'/'
```

BERT\_X KorQuAD-beginner NarrativeKoGPT2

**필요 패키지들 설치**

```
In [0]: !pip install -r drive/'My Drive'/'Colab Notebooks'/'NarrativeKoGPT2/requirements.txt
```

Collecting gluonnlp>=0.8.3

Downloading <https://files.pythonhosted.org/packages/c6/27/07b57d22496ed6c98b247e578712122402487f5c265ec70a747900f97060/gluonnlp-0.9.1.tar.gz> (252kB)

|██| 256kB 2.7MB/s

Collecting mxnet

Downloading <https://files.pythonhosted.org/packages/81/f5/d79b5b40735086ff1100c680703e0f3efc830fa455e268e9e96f3c857e93/mxnet-1.6.0-py2.py3-none-any.whl> (68.7MB)

|██| 68.7MB 44kB/s

Collecting sentencepiece>=0.1.6

Downloading [https://files.pythonhosted.org/packages/74/f4/2d5214cbf13d06e7cb2c20d84115ca25b53ea76fa1f0ade0e3c9749de214/sentencepiece-0.1.85-cp36-cp36m-manylinux1\\_x86\\_64.whl](https://files.pythonhosted.org/packages/74/f4/2d5214cbf13d06e7cb2c20d84115ca25b53ea76fa1f0ade0e3c9749de214/sentencepiece-0.1.85-cp36-cp36m-manylinux1_x86_64.whl) (1.0MB)

|██| 1.0MB 42.6MB/s

Requirement already satisfied: torch>=1.4.0 in /usr/local/lib/python3.6/dist-packages (from -r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 4)) (1.4.0)

Collecting transformers>=2.1.1

Downloading <https://files.pythonhosted.org/packages/13/33/ffb67897a6985a7b7d8e5e7878c3628678f553634bd3836404fef06ef19b/transformers-2.5.1-py3-none-any.whl> (499kB)

|██| 501kB 52.7MB/s

Requirement already satisfied: tqdm in /usr/local/lib/python3.6/dist-packages (from -r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 6)) (4.38.0)

Requirement already satisfied: numpy>=1.16.0 in /usr/local/lib/python3.6/dist-packages (from gluonnlp>=0.8.3->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 1)) (1.18.2)

Requirement already satisfied: cython in /usr/local/lib/python3.6/dist-packages (from gluonnlp>=0.8.3->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 1)) (0.29.15)

Requirement already satisfied: packaging in /usr/local/lib/python3.6/dist-packages (from gluonnlp>=0.8.3->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 1)) (20.3)

Collecting graphviz<0.9.0,>=0.8.1

Downloading <https://files.pythonhosted.org/packages/53/39/4ab213673844e0c004bed8a0781a0721a3f6bb23eb8854ee75c236428892/graphviz-0.8.4-py2.py3-none-any.whl>

Requirement already satisfied: requests<3,>=2.20.0 in /usr/local/lib/python3.6/dist-packages (from mxnet->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 2)) (2.21.0)

Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.6/dist-packages (from transformers>=2.1.1->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 5)) (2019.12.20)

Collecting sacremoses

Downloading <https://files.pythonhosted.org/packages/a6/b4/7a41d630547a4afd58143597d5a49e07bfd4c42914d8335b2a5657efc14b/sacremoses-0.0.38.tar.gz> (860kB)

|██| 870kB 43.9MB/s

Requirement already satisfied: boto3 in /usr/local/lib/python3.6/dist-packages (from transformers>=2.1.1->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 5)) (1.12.23)

Requirement already satisfied: filelock in /usr/local/lib/python3.6/dist-packages (from transformers>=2.1.1->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 5)) (3.0.12)

Collecting tokenizers==0.5.2

Downloading [https://files.pythonhosted.org/packages/d1/3f/73c881ea4723e43c1e9acf317cf407fab3a278daab3a69c98dcac511c04f/tokenizers-0.5.2-cp36-cp36m-manylinux1\\_x86\\_64.whl](https://files.pythonhosted.org/packages/d1/3f/73c881ea4723e43c1e9acf317cf407fab3a278daab3a69c98dcac511c04f/tokenizers-0.5.2-cp36-cp36m-manylinux1_x86_64.whl) (3.7MB)

|██| 3.7MB 37.1MB/s

Requirement already satisfied: pyparsing>=2.0.2 in /usr/local/lib/python3.6/dist-packages (from packaging->gluonnlp>=0.8.3->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 1)) (2.4.6)

Requirement already satisfied: six in /usr/local/lib/python3.6/dist-packages (from packaging->gluonnlp>=0.8.3->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 1)) (1.12.0)

Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.6/dist-packages (from requests<3,>=2.20.0->mxnet->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 2)) (2019.11.28)

Requirement already satisfied: chardet<3.1.0,>=3.0.2 in /usr/local/lib/python3.6/dist-packages (from requests<3,>=2.20.0->mxnet->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 2)) (3.0.4)

Requirement already satisfied: urllib3<1.25,>=1.21.1 in /usr/local/lib/python3.6/dist-packages (from requests<3,>=2.20.0->mxnet->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 2)) (1.24.3)

Requirement already satisfied: idna<2.9,>=2.5 in /usr/local/lib/python3.6/dist-packages (from requests<3,>=2.20.0->mxnet->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 2)) (2.8)

Requirement already satisfied: click in /usr/local/lib/python3.6/dist-packages (from sacremoses->transformers>=2.1.1->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 5)) (7.1.1)

Requirement already satisfied: joblib in /usr/local/lib/python3.6/dist-packages (from sacremoses->transformers>=2.1.1->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 5)) (0.14.1)

Requirement already satisfied: s3transfer<0.4.0,>=0.3.0 in /usr/local/lib/python3.6/dist-packages (from boto3->transformers>=2.1.1->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 5)) (0.3.3)

Requirement already satisfied: jmespath<1.0.0,>=0.7.1 in /usr/local/lib/python3.6/dist-packages (from boto3->transformers>=2.1.1->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 5)) (0.9.5)

Requirement already satisfied: botocore<1.16.0,>=1.15.23 in /usr/local/lib/python3.6/dist-packages (from boto3->transformers>=2.1.1->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 5)) (1.15.23)

Requirement already satisfied: python-dateutil<3.0.0,>=2.1 in /usr/local/lib/python3.6/dist-packages (from botocore<1.16.0,>=1.15.23->boto3->transformers>=2.1.1->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 5)) (2.8.1)

Requirement already satisfied: docutils<0.16,>=0.10 in /usr/local/lib/python3.6/dist-packages (from botocore<1.16.0,>=1.15.23->boto3->transformers>=2.1.1->-r drive/My Drive/Colab Notebooks/NarrativeKoGPT2/requirements.txt (line 5)) (0.15.2)

Building wheels for collected packages: gluonnlp, sacremoses

Building wheel for gluonnlp (setup.py) ... done

Created wheel for gluonnlp: filename=gluonnlp-0.9.1-cp36-cp36m-linux\_x86\_64.whl size=470973 sha256=ce0a0e8bff9f39914afe70e85415bf087a6dfdc32ee2642f8e8623e0d2847c2d

Stored in directory: /root/.cache/pip/wheels/af/60/16/1f8a40e68b85bd9bd7960e91830bca5e40cd113f3220b7e231

Building wheel for sacremoses (setup.py) ... done

Created wheel for sacremoses: filename=sacremoses-0.0.38-cp36-none-any.whl size=884628 sha256=b304eeacb4ed3c89b77e1e08a51c7e5176adbcfd8c0bb1a7523d252ae36192f

Stored in directory: /root/.cache/pip/wheels/6d/ec/1a/21b8912e35e02741306f35f66c785f3afe94de754a0eaf1422

Successfully built gluonnlp sacremoses

Installing collected packages: gluonnlp, graphviz, mxnet, sentencepiece, sacremoses, tokenizers, transformers

Found existing installation: graphviz 0.10.1

Uninstalling graphviz-0.10.1:

Successfully uninstalled graphviz-0.10.1

Successfully installed gluonnlp-0.9.1 graphviz-0.8.4 mxnet-1.6.0 sacremoses-0.0.38 sentencepiece-0.1.85 tokenizers-0.5.2 transformers-2.5.1

## 시스템 경로 추가

```
In [0]: import os
import sys
sys.path.append('drive/My Drive/Colab Notebooks/')
```

## 2.KoGPT2 Text Generation

### 2.1.Import Package

```
In [0]: import random
import torch
from torch.utils.data import DataLoader # 데이터로더
from gluonnlp.data import SentencepieceTokenizer
from NarrativeKoGPT2.kogpt2.utils import get_tokenizer
from NarrativeKoGPT2.kogpt2.utils import download, tokenizer
from NarrativeKoGPT2.model.torch_gpt2 import GPT2Config, GPT2LMHeadModel
from NarrativeKoGPT2.util.data import NovelDataset
import gluonnlp
```

The default version of TensorFlow in Colab will soon switch to TensorFlow 2.x.

We recommend you [upgrade now](#) or ensure your notebook will continue to use TensorFlow 1.x via the `%tensorflow_version 1.x` magic: [more info](#).

### 2.2. koGPT-2 Config

```
In [0]: ctx= 'cpu' #'cuda' #'cpu' #학습 Device CPU or GPU. colab의 경우 GPU 사용
cachedir='~/kogpt2/' # KoGPT-2 모델 다운로드 경로
epoch =200 # 학습 epoch
save_path = 'drive/My Drive/Colab Notebooks/NarrativeKoGPT2/checkpoint/'
load_path = 'drive/My Drive/Colab Notebooks/NarrativeKoGPT2/checkpoint/narrativeKoG
#use_cuda = True # Colab내 GPU 사용을 위한 값

pytorch_kogpt2 = {
    'url':
    'https://kobert.blob.core.windows.net/models/kogpt2/pytorch/pytorch_kogpt2_676e
    'fname': 'pytorch_kogpt2_676e9bcfa7.params',
    'chksum': '676e9bcfa7'
}
kogpt2_config = {
    "initializer_range": 0.02,
    "layer_norm_epsilon": 1e-05,
    "n_ctx": 1024,
    "n_embd": 768,
    "n_head": 12,
    "n_layer": 12,
    "n_positions": 1024,
```

```
"vocab_size": 50000
}
```

## 2.3 Model and Vocab Download

```
In [0]: # download model
model_info = pytorch_kogpt2
model_path = download(model_info['url'],
                      model_info['fname'],
                      model_info['chksum'],
                      cachedir=cachedir)

# download vocab
vocab_info = tokenizer
vocab_path = download(vocab_info['url'],
                    vocab_info['fname'],
                    vocab_info['chksum'],
                    cachedir=cachedir)
```



```
[REDACTED]
[REDACTED]
```

## 2.4.KoGPT-2 Model Vocab

추론 및 학습 재개를 위한 모델 불러오기 저장하기

```
torch.save({
    'epoch': epoch,
    'model_state_dict': model.state_dict(),
    'optimizer_state_dict': optimizer.state_dict(),
    'loss': loss,
    ...
}, PATH)
```

불러오기

```
model = TheModelClass(*args, **kwargs)
optimizer = TheOptimizerClass(*args, **kwargs)

checkpoint = torch.load(PATH)
model.load_state_dict(checkpoint['model_state_dict'])
optimizer.load_state_dict(checkpoint['optimizer_state_dict'])
epoch = checkpoint['epoch']
loss = checkpoint['loss']

model.eval()
# - or -
model.train()
```

```
In [0]: # Device 설정
device = torch.device(ctx)
# 저장한 Checkpoint 불러오기
checkpoint = torch.load(load_path, map_location=device)
```

```
# KoGPT-2 언어 모델 학습을 위한 GPT2LMHeadModel 선언
kogpt2model = GPT2LMHeadModel(config=GPT2Config.from_dict(kogpt2_config))
kogpt2model.load_state_dict(checkpoint['model_state_dict'])

kogpt2model.eval()
vocab_b_obj = gluonnlp.vocab.BERTVocab.from_sentencepiece(vocab_path,
                                                         mask_token=None,
                                                         sep_token=None,
                                                         cls_token=None,
                                                         unknown_token='<unk>',
                                                         padding_token='<pad>',
                                                         bos_token='<s>',
                                                         eos_token='</s>')
```

## 2.5. Tokenizer

```
In [0]: tok_path = get_tokenizer()
        model, vocab = kogpt2model, vocab_b_obj
        tok = SentencepieceTokenizer(tok_path)
```

using cached model

## 2.6. NarrativeKoGPT-2 Text Generation

```
In [0]: sent = input('문장 입력: ')

        toked = tok(sent)
        count = 0
        output_size = 200 # 출력하고자 하는 토큰 갯수

        file 1:
        input_ids = torch.tensor([vocab[vocab.bos_token],] + vocab[toked]).unsqueeze(0)
        predicts = model(input_ids)
        pred = predicts[0]

        last_pred = pred.squeeze()[-1]
        # top_p 샘플링 방법
        # sampling.py를 통해 random, top-k, top-p 선택 가능.
        gen = sampling.top_p(last_pred, vocab, 0.9)
        # gen = sampling.top_k(last_pred, vocab, 5)

        if count > output_size:
            sent += gen.replace('_', ' ')
            toked = tok(sent)
            count = 0
            break
        sent += gen.replace('_', ' ')
        toked = tok(sent)
        count += 1

        for s in kss.split_sentences(sent):
            print(s)
```

