# Artificial Intelligence for Text Analytics

# Text Classification and Sentiment Analysis

**Min-Yuh Day, Ph.D,**
**Associate Professor**

**Institute of Information Management**, **National Taipei University**

https://web.ntpu.edu.tw/~myday

https://meet.google.com/
paj-zhhj-mya

2022-03-29

# Syllabus

Week    Date    Subject/Topics

1   2022/02/22   Introduction to Artificial Intelligence for Text Analytics

2   2022/03/01   Foundations of Text Analytics:
                 Natural Language Processing (NLP)

3   2022/03/08   Python for Natural Language Processing

4   2022/03/15   Natural Language Processing with Transformers

5   2022/03/22   Case Study on Artificial Intelligence for Text Analytics I

6   2022/03/29   Text Classification and Sentiment Analysis

# Syllabus

Week    Date    Subject/Topics

7   2022/04/05   Tomb-Sweeping Day (Holiday, No Classes)

8   2022/04/12   Midterm Project Report

9   2022/04/19   Multilingual Named Entity Recognition (NER),
                 Text Similarity and Clustering

10   2022/04/26   Text Summarization and Topic Models

11   2022/05/03   Text Generation

12   2022/05/10   Case Study on Artificial Intelligence for Text Analytics II

# Syllabus

Week    Date    Subject/Topics

13   2022/05/17   Question Answering and Dialogue Systems

14   2022/05/24   Deep Learning, Transfer Learning,
                  Zero-Shot, and Few-Shot Learning for Text Analytics

15   2022/05/31   Final Project Report I

16   2022/06/07   Final Project Report II

17   2022/06/14   Self-learning

18   2022/06/21   Self-learning

# **Text Classification and Sentiment Analysis**
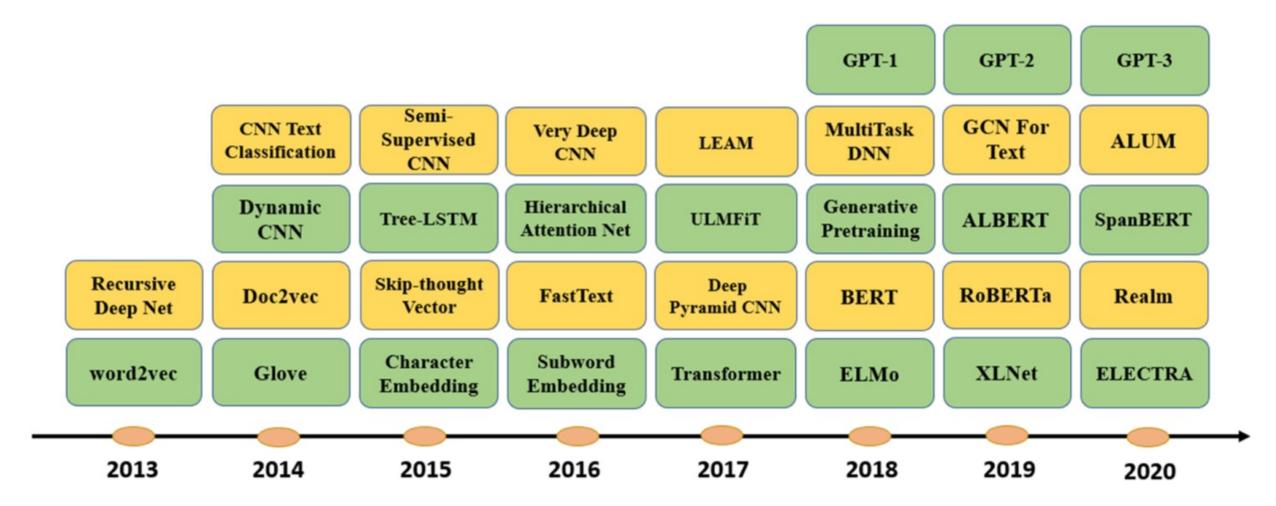
# Outline

- **Text Classification and Sentiment Analysis**
  - **Dataset**
  - **Tokenizer**
  - **Training a Text Classifier**
  - **Fine-Tuning Transformers**

# Text Classification (TC) Tasks

- **Sentiment Analysis**

- **News Categorization**

- **Product Categorization**

- **Topic Analysis**

  - Topic Classification: "customer support" or "ease of use"

- **Natural language inference (NLI)**

  - recognizing textual entailment (RTE)

  - entailment, contradiction, and neutral

# Deep learning models for text embedding and classification

# Text Classification Models on Sentiment Analysis

| Method | IMDB | SST-2 | Amazon-2 | Amazon-5 | Yelp-2 | Yelp-5 |
|---|---|---|---|---|---|---|
| *Naive Bayes* [43] | - | 81.80 | - | - | - | - |
| *LDA* [214] | 67.40 | - | - | - | - | - |
| *BoW+SVM* [31] | 87.80 | - | - | - | - | - |
| *tf.Δ idf* [215] | 88.10 | - | - | - | - | - |
| Char-level CNN [50] | - | - | 94.49 | 59.46 | 95.12 | 62.05 |
| Deep Pyramid CNN [49] | - | 84.46 | 96.68 | 65.82 | 97.36 | 69.40 |
| ULMFiT [216] | 95.40 | - | - | - | 97.84 | 70.02 |
| BLSTM-2DCNN [40] | - | 89.50 | - | - | - | - |
| Neural Semantic Encoder [95] | - | 89.70 | - | - | - | - |
| BCN+Char+CoVe [217] | 91.80 | 90.30 | - | - | - | - |
| GLUE ELMo baseline [22] | - | 90.40 | - | - | - | - |
| BERT ELMo baseline [7] | - | 90.40 | - | - | - | - |
| CCCapsNet [76] | - | - | 94.96 | 60.95 | 96.48 | 65.85 |
| Virtual adversarial training [173] | 94.10 | - | - | - | - | - |
| Block-sparse LSTM [218] | 94.99 | 93.20 | - | - | 96.73 | |
| BERT-base [7, 154] | 95.63 | 93.50 | 96.04 | 61.60 | 98.08 | 70.58 |
| BERT-large [7, 154] | 95.79 | 94.9 | 96.07 | 62.20 | 98.19 | 71.38 |
| ALBERT [147] | - | 95.20 | - | - | - | - |
| Multi-Task DNN [23] | 83.20 | 95.60 | - | - | - | - |
| Snorkel MeTaL [219] | - | 96.20 | - | - | - | - |
| BERT Finetune + UDA [220] | 95.80 | | 96.50 | 62.88 | 97.95 | 62.92 |
| RoBERTa (+additional data) [146] | - | 96.40 | - | - | - | - |
| XLNet-Large (ensemble) [156] | 96.21 | 96.80 | 97.60 | 67.74 | 98.45 | 72.20 |

# Classification Models on News Categorization, and Topic Classification

| Method | News Categorization | | | Topic Classification | |
|---|---|---|---|---|---|
| | AG News | 20NEWS | Sogou News | DBpedia | Ohsumed |
| *Hierarchical Log-bilinear Model* [221] | - | - | - | - | 52 |
| Text GCN [107] | 67.61 | 86.34 | - | - | 68.36 |
| Simplfied GCN [108] | - | 88.50 | - | - | 68.50 |
| Char-level CNN [50] | 90.49 | - | 95.12 | 98.45 | - |
| CCCapsNet [76] | 92.39 | - | 97.25 | 98.72 | - |
| LEAM [84] | 92.45 | 81.91 | - | 99.02 | 58.58 |
| fastText [30] | 92.50 | - | 96.80 | 98.60 | 55.70 |
| CapsuleNet B [71] | 92.60 | - | - | - | - |
| Deep Pyramid CNN [49] | 93.13 | - | 98.16 | 99.12 | - |
| ULMFiT [216] | 94.99 | - | - | 99.20 | - |
| L MIXED [174] | 95.05 | - | - | 99.30 | - |
| BERT-large [220] | - | - | - | 99.32 | - |
| XLNet [156] | 95.51 | - | - | 99.38 | - |

# Classification Models on Natural Language Inference (NLI)

| Method | SNLI Accuracy | MultiNLI Matched | MultiNLI Mismatched |
|---|---|---|---|
| *Unigrams Features* [208] | 71.6 | — | — |
| *Lexicalized* [208] | 78.2 | — | — |
| LSTM encoders (100D) [208] | 77.6 | — | — |
| Tree-based CNN [61] | 82.1 | — | — |
| biLSTM Encoder [209] | 81.5 | 67.5 | 67.1 |
| Neural Semantic Encoders (300D) [95] | 84.6 | — | — |
| RNN-based Sentence Encoder [224] | 85.5 | 73.2 | 73.6 |
| DiSAN (300D) [81] | 85.6 | — | — |
| Decomposable Attention Model [92] | 86.3 | — | — |
| Reinforced Self-Attention (300D) [177] | 86.3 | — | — |
| Generalized Pooling (600D) [93] | 86.6 | 73.8 | 74.0 |
| Bilateral multi-perspective matching [41] | 87.5 | — | — |
| Multiway Attention Network [87] | 88.3 | 78.5 | 77.7 |
| ESIM + ELMo [4] | 88.7 | 72.9 | 73.4 |
| DMAN with Reinforcement Learning [225] | 88.8 | 88.8 | 78.9 |
| BiLSTM + ELMo + Attn [22] | — | 74.1 | 74.5 |
| Fine-Tuned LM-Pretrained Transformer [6] | 89.9 | 82.1 | 81.4 |
| Multi-Task DNN [23] | 91.6 | 86.7 | 86.0 |
| SemBERT [155] | 91.9 | 84.4 | 84.0 |
| RoBERTa [146] | 92.6 | 90.8 | 90.2 |
| XLNet [156] | — | 90.2 | 89.8 |

# General Language Understanding Evaluation (GLUE) benchmark
# GLUE Test results

| System | MNLI-(m/mm) | QQP | QNLI | SST-2 | CoLA | STS-B | MRPC | RTE | Average |
|---|---|---|---|---|---|---|---|---|---|
| | 392k | 363k | 108k | 67k | 8.5k | 5.7k | 3.5k | 2.5k | - |
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.9 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 88.1 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.2 |
| BERT$_{BASE}$ | 84.6/83.4 | 71.2 | 90.1 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| BERT$_{LARGE}$ | **86.7/85.9** | **72.1** | **91.1** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **81.9** |

**MNLI**: Multi-Genre Natural Language Inference

**QQP**: Quora Question Pairs

**QNLI**: Question Natural Language Inference

**SST-2**: The Stanford Sentiment Treebank

**CoLA**: The Corpus of Linguistic Acceptability

**STS-B**:The Semantic Textual Similarity Benchmark

**MRPC**: Microsoft Research Paraphrase Corpus

**RTE**: Recognizing Textual Entailment

# Emotions

| | |
|---|---|
| Love | Anger |
| Joy | Sadness |
| Surprise | Fear |

# Example of Opinion:
# review segment on iPhone

"I bought an iPhone a few days ago.

It was such a nice phone.

The touch screen was really cool.

The voice quality was clear too.

However, my mother was mad with me as I did not tell her before I bought it.

She also thought the phone was too expensive, and wanted me to return it to the shop. … "

# Example of Opinion: review segment on iPhone

"(1) I bought an <u>iPhone</u> a few days ago.

(2) **It was such a nice phone.**

(3) **The <u>touch screen</u> was really cool.**

**+Positive Opinion**

(4) **The <u>voice quality</u> was clear too.**

(5) However, my mother was mad with me as I did not tell her before I bought it.

(6) She also thought the phone was too <u>expensive</u>, and wanted me to return it to the shop. … "

**-Negative Opinion**

# Sentiment Analysis



Subjectivity Classification

Sentiment Classification

Review Usefulness Measurement

Opinion Spam Detection

Lexicon Creation

Aspect Extraction

Sentiment Analysis

Application

Polarity Determination

Vagueness resolution in opinionated text

Multi- & Cross-Lingual SC

Cross-domain SC

**Tasks**

**Approaches**

Machine Learning based

Lexicon based

Hybrid approaches

Ontology based

Non-Ontology based

# Sentiment Classification Techniques

17

# P–N Polarity and S–O Polarity Relationship



Source: Ramesh Sharda, Dursun Delen, and Efraim Turban (2017), Business Intelligence, Analytics, and Data Science: A Managerial Perspective, 4th Edition, Pearson
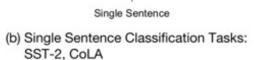
# Lewis Tunstall, Leandro von Werra, and Thomas Wolf (2022),
# Natural Language Processing with Transformers:
## Building Language Applications with Hugging Face,
### O'Reilly Media.

# ULMFiT: 3 Steps
# Transfer Learning in NLP



**1. Pretraining**          **2. Domain adaptation**          **3. Fine-tuning**

# An overview of the Hugging Face Ecosystem

# A typical pipeline for training transformer models

with the  Datasets,  Tokenizers, and  Transformers libraries

| Datasets | Tokenizers | Transformers | Datasets |
|---|---|---|---|
| **Load and process datasets** | **Tokenize input texts** | **Load models, train and infer** | **Load metrics evaluate models** |

# NLP with Transformers

```
!git clone https://github.com/nlp-with-transformers/notebooks.git
%cd notebooks
from install import *
install_requirements()
```

```
from utils import *
setup_chapter()
```

# Text Classification

```
text = """"Dear Amazon, last week I ordered an Optimus Prime action figure \
from your online store in Germany. Unfortunately, when I opened the package, \
I discovered to my horror that I had been sent an action figure of Megatron \
instead! As a lifelong enemy of the Decepticons, I hope you can understand my \
dilemma. To resolve the issue, I demand an exchange of Megatron for the \
Optimus Prime figure I ordered. Enclosed are copies of my records concerning \
this purchase. I expect to hear from you soon. Sincerely, Bumblebee."""
```

# Text Classification

```
text = """Dear Amazon, last week I ordered an Optimus Prime action figure \
from your online store in Germany. Unfortunately, when I opened the package, \
I discovered to my horror that I had been sent an action figure of Megatron \
instead! As a lifelong enemy of the Decepticons, I hope you can understand my \
dilemma. To resolve the issue, I demand an exchange of Megatron for the \
Optimus Prime figure I ordered. Enclosed are copies of my records concerning \
this purchase. I expect to hear from you soon. Sincerely, Bumblebee."""
```

```
from transformers import pipeline
classifier = pipeline("text-classification")
```

```
import pandas as pd
outputs = classifier(text)
pd.DataFrame(outputs)
```

|   | label | score |
|---|-------|-------|
| 0 | NEGATIVE | 0.901546 |

# Text Classification

```
from transformers import pipeline
classifier = pipeline("text-classification")
```

```
import pandas as pd
outputs = classifier(text)
pd.DataFrame(outputs)
```

|   | label | score |
|---|-------|-------|
| 0 | NEGATIVE | 0.901546 |

# Fine-tuning BERT on NLP Tasks



(a) Sentence Pair Classification Tasks: MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG

(b) Single Sentence Classification Tasks: SST-2, CoLA

(c) Question Answering Tasks: SQuAD v1.1

(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

# BERT Sequence-level tasks



(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG

(b) Single Sentence Classification Tasks:
SST-2, CoLA

# BERT Token-level tasks



(c) Question Answering Tasks: SQuAD v1.1

(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

# Sentiment Analysis:
# Single Sentence Classification



(b) Single Sentence Classification Tasks:
SST-2, CoLA

# Character Tokenization

```python
text = "Tokenizing text is a core task of NLP."
tokenized_text = list(text)
print(tokenized_text)
```

```
['T', 'o', 'k', 'e', 'n', 'i', 'z', 'i', 'n', 'g', ' ', 't', 'e', 'x', 't', ' ',
 'i', 's', ' ', 'a', ' ', 'c', 'o', 'r', 'e', ' ', 't', 'a', 's', 'k', ' ', 'o',
 'f', ' ', 'N', 'L', 'P', '.']
```

```python
token2idx = {ch: idx for idx, ch in enumerate(sorted(set(tokenized_text)))}
print(token2idx)
```

```
{' ': 0, '.': 1, 'L': 2, 'N': 3, 'P': 4, 'T': 5, 'a': 6, 'c': 7, 'e': 8, 'f': 9,
 'g': 10, 'i': 11, 'k': 12, 'n': 13, 'o': 14, 'r': 15, 's': 16, 't': 17, 'x': 18,
 'z': 19}
```

```python
input_ids = [token2idx[token] for token in tokenized_text]
print(input_ids)
```

```
[5, 14, 12, 8, 13, 11, 19, 11, 13, 10, 0, 17, 8, 18, 17, 0, 11, 16, 0, 6, 0, 7,
14, 15, 8, 0, 17, 6, 16, 12, 0, 14, 9, 0, 3, 2, 4, 1]
```

# Word Tokenization

```python
text = "Tokenizing text is a core task of NLP."
tokenized_text = text.split()
print(tokenized_text)
```

```
['Tokenizing', 'text', 'is', 'a', 'core', 'task', 'of', 'NLP.']
```

# Subword Tokenization

```python
from transformers import AutoTokenizer
model_ckpt = "distilbert-base-uncased"
tokenizer = AutoTokenizer.from_pretrained(model_ckpt)
```

```python
text = "Tokenizing text is a core task of NLP."
encoded_text = tokenizer(text)
print(encoded_text)
```

```
{'input_ids': [101, 19204, 6026, 3793, 2003, 1037, 4563, 4708, 1997, 17953, 2361, 1012, 102], 'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]}
```

```python
tokens = tokenizer.convert_ids_to_tokens(encoded_text.input_ids)
print(tokens)
```

```
['[CLS]', 'token', '##izing', 'text', 'is', 'a', 'core', 'task', 'of', 'nl', '##p', '.', '[SEP]']
```
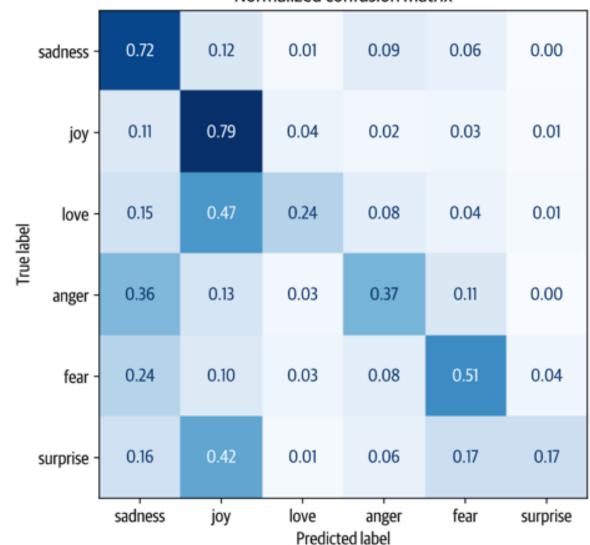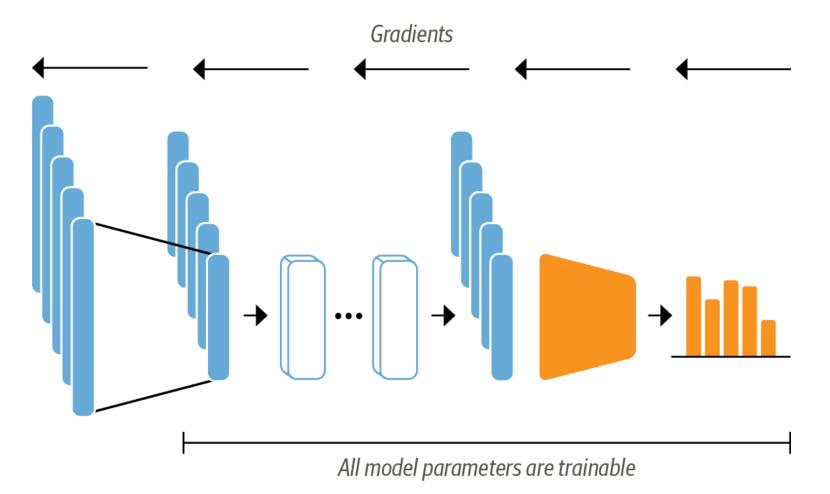
Source: Lewis Tunstall, Leandro von Werra, and Thomas Wolf (2022), Natural Language Processing with Transformers:  Building Language Applications with Hugging Face,  O'Reilly Media.

https://github.com/nlp-with-transformers/notebooks

33

# Subword Tokenization

```
print(tokenizer.convert_tokens_to_string(tokens))
```

```
[CLS] tokenizing text is a core task of nlp. [SEP]
```

```
tokenizer.vocab_size
```

```
30522
```

```
tokenizer.model_max_length
```

```
512
```

# Tokenizing the Whole Dataset

```python
def tokenize(batch):
    return tokenizer(batch["text"], padding=True, truncation=True)
```

```python
print(tokenize(emotions["train"][:2]))
```

```
{'input_ids': [[101, 1045, 2134, 2102, 2514, 26608, 102, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [101, 1045, 2064, 2175, 2013, 3110, 2061,
20625, 2000, 2061, 9636, 17772, 2074, 2013, 2108, 2105, 2619, 2040, 14977,
1998, 2003, 8300, 102]], 'attention_mask': [[1, 1, 1, 1, 1, 1, 1, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]]}
```

```python
tokens2ids = list(zip(tokenizer.all_special_tokens,
tokenizer.all_special_ids))
data = sorted(tokens2ids, key=lambda x : x[-1])
df = pd.DataFrame(data, columns=["Special Token", "Special Token ID"])
df.T
```

# From Text to Tokens

**For each batch, the input sequences are padded to the maximum sequence length in the batch; the attention mask is used in the model to ignore the padded areas of the input tensors**

# Training a Text Classifier

**The architecture used for sequence classification with an encoder-based transformer ; it consists of the model's pretrained body combined with a custom classification head**



| Token encodings | Token embeds | Encoder stack | Hidden states | Classification head | Predictions |

# Transformers as Feature Extractors

**In the feature-based approach, the DistilBERT model is frozen and just provides features for a classifier**

# Training a Simple Classifier



Normalized confusion matrix

# Fine-Tuning Transformers

## When using the fine-tuning approach the whole DistilBERT model is trained along with the classification head

https://github.com/nlp-with-transformers/notebooks

# Fine-Tuning Transformers
# Loading a pretrained model

```python
from transformers import AutoModelForSequenceClassification

num_labels = 6
model = (AutoModelForSequenceClassification
         .from_pretrained(model_ckpt, num_labels=num_labels)
         .to(device))
```

# Defining the performance metrics

```python
from sklearn.metrics import accuracy_score, f1_score

def compute_metrics(pred):
    labels = pred.label_ids
    preds = pred.predictions.argmax(-1)
    f1 = f1_score(labels, preds, average="weighted")
    acc = accuracy_score(labels, preds)
    return {"accuracy": acc, "f1": f1}
```

# Train the model

```
from huggingface_hub import notebook_login

notebook_login()
```

# Train the model

```python
from transformers import Trainer, TrainingArguments

batch_size = 64
logging_steps = len(emotions_encoded["train"]) // batch_size
model_name = f"{model_ckpt}-finetuned-emotion"
training_args = TrainingArguments(output_dir=model_name,
                    num_train_epochs=2,
                    learning_rate=2e-5,
                    per_device_train_batch_size=batch_size,
                    per_device_eval_batch_size=batch_size,
                    weight_decay=0.01,
                    evaluation_strategy="epoch",
                    disable_tqdm=False,
                    logging_steps=logging_steps,
                    push_to_hub=True,
                    log_level="error")
```

# Train the model

```python
from transformers import Trainer

trainer = Trainer(model=model, args=training_args,
        compute_metrics=compute_metrics,
        train_dataset=emotions_encoded["train"],
        eval_dataset=emotions_encoded["validation"],
        tokenizer=tokenizer)
trainer.train();
```

[500/500 01:48, Epoch 2/2]

| Epoch | Training Loss | Validation Loss | Accuracy | F1 |
|---|---|---|---|---|
| 1 | 0.840900 | 0.327445 | 0.896500 | 0.892285 |
| 2 | 0.255000 | 0.220472 | 0.922500 | 0.922550 |

# Train the model

```
preds_output =
trainer.predict(emotions_encoded["validation"])
```

```
preds_output.metrics
```

{'test_loss': 0.22047173976898193, 'test_accuracy': 0.9225, 'test_f1': 0.9225500751072866, 'test_runtime': 1.6357, 'test_samples_per_second': 1222.725, 'test_steps_per_second': 19.564}

```
y_preds = np.argmax(preds_output.predictions, axis=1)
```

```
plot_confusion_matrix(y_preds, y_valid, labels)
```

# Fine-Tuning Transformers



Normalized confusion matrix

# A Visual Guide to
# Using BERT for the First Time
**(Jay Alammar, 2019)**

# Sentiment Classification: SST2
# Sentences from movie reviews

| sentence | label |
|---|---|
| a stirring , funny and finally transporting re imagining of beauty and the beast and 1930s horror films | 1 |
| apparently reassembled from the cutting room floor of any given daytime soap | 0 |
| they presume their audience won't sit still for a sociology lesson | 0 |
| this is a visually stunning rumination on love , memory , history and the war between art and commerce | 1 |
| jonathan parker 's bartleby should have been the be all end all of the modern office anomie films | 1 |

# Movie Review Sentiment Classifier

# Movie Review Sentiment Classifier

# Movie Review Sentiment Classifier Model Training



Movie Review Sentiment Classifier

DistilBERT

Already (pre-)trained

Logistic Regression

We will train in this tutorial

# Step # 1 Use distilBERT to Generate Sentence Embeddings

# Step #2:Test/Train Split for Model #2, Logistic Regression



Step #2: Test/Train Split for model #2, logistic regression

# Step #3 Train the logistic regression model using the training set

# Tokenization

[CLS] a visually stunning rum ##ination on love [SEP]
a visually stunning rumination on love



**Tokenization**
DistilBertTokenizer

| [CLS] | a | visually | stunning | rum | ##ination | on | love | [SEP] |

2) Add [CLS] and [SEP] tokens

| a | visually | stunning | rum | ##ination | on | love |

1) Break words into tokens

**Tokenize**

"a visually stunning rumination on love"

# Tokenization

```
tokenizer.encode("a visually stunning rumination on love",
                 add_special_tokens=True)
```



**Tokenization**
DistilBertTokenizer

| 101 | 1037 | 17453 | 14726 | 19379 | 12758 | 2006 | 2293 | 102 |

3) substitute tokens with their ids

| [CLS] | a | visually | stunning | rum | ##ination | on | love | [SEP] |

2) Add [CLS] and [SEP] tokens

| a | visually | stunning | rum | ##ination | on | love |

1) Break words into tokens

Tokenize

"a visually stunning rumination on love"

# Tokenization for BERT Model

# Flowing Through DistilBERT (768 features)

59

# Model #1 Output **Class** vector as Model #2 Input

# Fine-tuning BERT on Single Sentence Classification Tasks

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# Model #1 Output **Class** vector as Model #2 Input

Source: Jay Alammar (2019), A Visual Guide to Using BERT for the First Time, http://jalammar.github.io/a-visual-guide-to-using-bert-for-the-first-time/
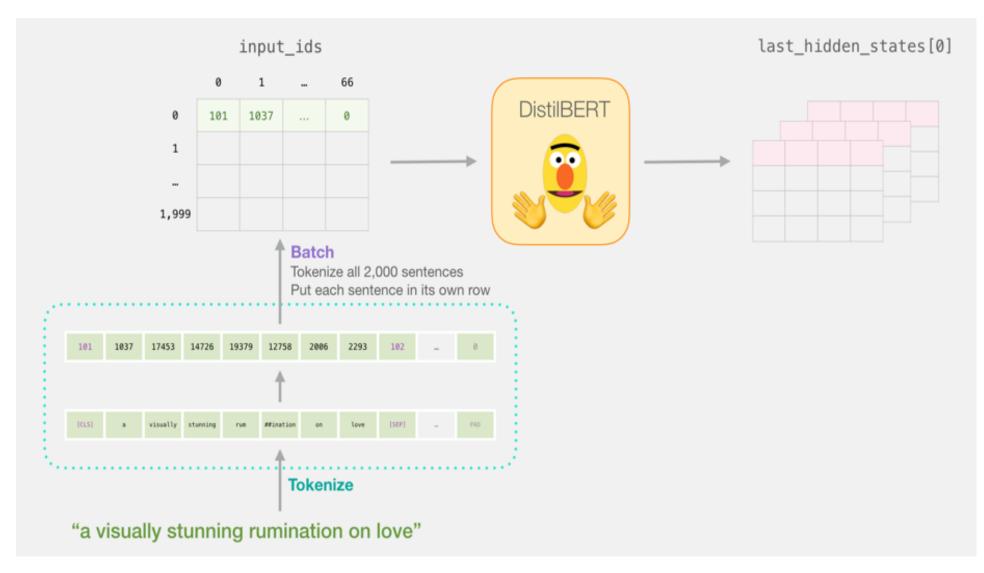
# Logistic Regression Model to classify **Class** vector

63

```
df = pd.read_csv('https://github.com/clairett/pytorch-
sentiment-classification/raw/master/data/SST2/train.tsv',
delimiter='\t', header=None)

df.head()
```

|   | 0 | 1 |
|---|---|---|
| 0 | a stirring , funny and finally transporting re... | 1 |
| 1 | apparently reassembled from the cutting room f... | 0 |
| 2 | they presume their audience wo n't sit still f... | 0 |
| 3 | this is a visually stunning rumination on love... | 1 |
| 4 | jonathan parker 's bartleby should have been t... | 1 |

# Tokenization

```
tokenized = df[0].apply((lambda x: tokenizer.encode(x,
add_special_tokens=True)))
```



Source: Jay Alammar (2019), A Visual Guide to Using BERT for the First Time,
http://jalammar.github.io/a-visual-guide-to-using-bert-for-the-first-time/

# BERT Input Tensor



BERT/DistilBERT Input Tensor

# Processing with DistilBERT

```
input_ids = torch.tensor(np.array(padded))
last_hidden_states = model(input_ids)
```

# Unpacking the BERT output tensor

# Sentence to last_hidden_state[0]

# BERT's output for the [CLS] tokens

```
# Slice the output for the first position for all the
sequences, take all hidden unit outputs
features = last_hidden_states[0][:,0,:].numpy()
```

# The tensor sliced from BERT's output
# Sentence Embeddings

# Dataset for Logistic Regression (768 Features)

**The features are the output vectors of BERT for the [CLS] token (position #0)**

```
labels = df[1]
train_features, test_features, train_labels, test_labels =
train_test_split(features, labels)
```

Step #2: Test/Train Split for model #2, logistic regression

# Score Benchmarks
# Logistic Regression Model
# on SST-2 Dataset

```
# Training
lr_clf = LogisticRegression()
lr_clf.fit(train_features, train_labels)

#Testing
lr_clf.score(test_features, test_labels)

# Accuracy: 81%
# Highest accuracy: 96.8%
# Fine-tuned DistilBERT: 90.7%
# Full size BERT model: 94.9%
```

# Sentiment Classification: SST2
# Sentences from movie reviews

| sentence | label |
|---|:---:|
| a stirring , funny and finally transporting re imagining of beauty and the beast and 1930s horror films | 1 |
| apparently reassembled from the cutting room floor of any given daytime soap | 0 |
| they presume their audience won't sit still for a sociology lesson | 0 |
| this is a visually stunning rumination on love , memory , history and the war between art and commerce | 1 |
| jonathan parker 's bartleby should have been the be all end all of the modern office anomie films | 1 |

# A Visual Notebook to Using BERT for the First Time

# Hugging Face Tasks
# Natural Language Processing



**Text Classification**

3345 models

**Token Classification**

1492 models

**Question Answering**

1140 models

**Translation**

1467 models

**Summarization**

323 models

**Text Generation**

3959 models

**Fill-Mask**

2453 models

**Sentence Similarity**

352 models

https://huggingface.co/tasks

# NLP with Transformers Github

# NLP with Transformers Github Notebooks



## Running on a cloud platform

To run these notebooks on a cloud platform, just click on one of the badges in the table below:

| Chapter | Colab | Kaggle | Gradient | Studio Lab |
| --- | --- | --- | --- | --- |
| Introduction | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |
| Text Classification | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |
| Transformer Anatomy | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |
| Multilingual Named Entity Recognition | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |
| Text Generation | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |
| Summarization | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |
| Question Answering | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |
| Making Transformers Efficient in Production | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |
| Dealing with Few to No Labels | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |
| Training Transformers from Scratch | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |
| Future Directions | Open in Colab | Open in Kaggle | Run on Gradient | Open Studio Lab |

Nowadays, the GPUs on Colab tend to be K80s (which have limited memory), so we recommend using Kaggle, Gradient, or SageMaker Studio Lab. These platforms tend to provide more performant GPUs like P100s, all for free!

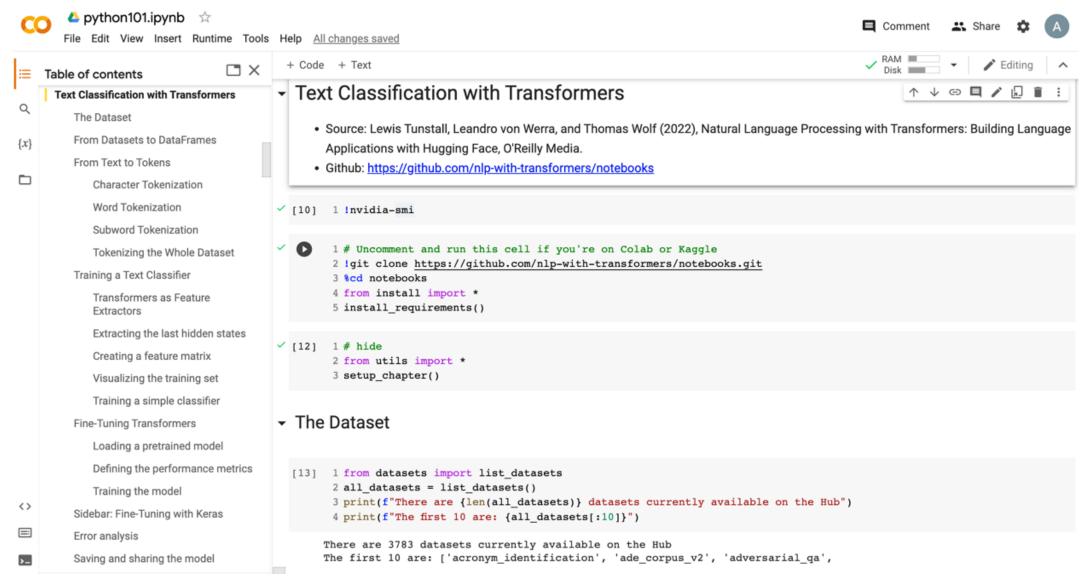https://github.com/nlp-with-transformers/notebooks

# Python in Google Colab (Python101)

# Python in Google Colab (Python101)

https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT

## Text Classification with Transformers

- Source: Lewis Tunstall, Leandro von Werra, and Thomas Wolf (2022), Natural Language Processing with Transformers: Building Language Applications with Hugging Face, O'Reilly Media.
- Github: https://github.com/nlp-with-transformers/notebooks

```
[10]  1 !nvidia-smi
```

```
      1 # Uncomment and run this cell if you're on Colab or Kaggle
      2 !git clone https://github.com/nlp-with-transformers/notebooks.git
      3 %cd notebooks
      4 from install import *
      5 install_requirements()
```

```
[12]  1 # hide
      2 from utils import *
      3 setup_chapter()
```

## The Dataset

```
[13]  1 from datasets import list_datasets
      2 all_datasets = list_datasets()
      3 print(f"There are {len(all_datasets)} datasets currently available on the Hub")
      4 print(f"The first 10 are: {all_datasets[:10]}")

There are 3783 datasets currently available on the Hub
The first 10 are: ['acronym_identification', 'ade_corpus_v2', 'adversarial_qa',
```

https://tinyurl.com/aintpupython101

81

# Summary

- **Text Classification and Sentiment Analysis**
  - **Dataset**
  - **Tokenizer**
  - **Training a Text Classifier**
  - **Fine-Tuning Transformers**

# References

- Lewis Tunstall, Leandro von Werra, and Thomas Wolf (2022), Natural Language Processing with Transformers:  Building Language Applications with Hugging Face, O'Reilly Media.
- Denis Rothman (2021), Transformers for Natural Language Processing: Build innovative deep neural network architectures for NLP with Python, PyTorch, TensorFlow, BERT, RoBERTa, and more, Packt Publishing.
- Savaş Yıldırım and Meysam Asgari-Chenaghlu (2021), Mastering Transformers: Build state-of-the-art models from scratch with advanced natural language processing techniques, Packt Publishing.
- Shervin Minaee, Nal Kalchbrenner, Erik Cambria, Narjes Nikzad, Meysam Chenaghlu, and Jianfeng Gao (2021), "Deep learning--based text classification: a comprehensive review." ACM Computing Surveys (CSUR) 54, no. 3 (2021): 1-40.
- Sowmya Vajjala, Bodhisattwa Majumder, Anuj Gupta (2020), Practical Natural Language Processing: A Comprehensive Guide to Building Real-World NLP Systems, O'Reilly Media.
- Ramesh Sharda, Dursun Delen, and Efraim Turban (2017), Business Intelligence, Analytics, and Data Science: A Managerial Perspective, 4th Edition, Pearson.
- Dipanjan Sarkar (2019), Text Analytics with Python: A Practitioner's Guide to Natural Language Processing, Second Edition. APress.
- Benjamin Bengfort, Rebecca Bilbro, and Tony Ojeda (2018), Applied Text Analysis with Python: Enabling Language-Aware Data Products with Machine Learning, O'Reilly.
- Charu C. Aggarwal (2018), Machine Learning for Text, Springer.
- Gabe Ignatow and Rada F. Mihalcea (2017), An Introduction to Text Mining: Research Design, Data Collection, and Analysis, SAGE Publications.
- Rajesh Arumugam (2018), Hands-On Natural Language Processing with Python: A practical guide to applying deep learning architectures to your NLP applications, Packt.
- Jake VanderPlas (2016), Python Data Science Handbook: Essential Tools for Working with Data, O'Reilly Media.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv preprint arXiv:1810.04805.
- The Super Duper NLP Repo, https://notebooks.quantumstat.com/
- Jay Alammar (2018), The Illustrated Transformer, http://jalammar.github.io/illustrated-transformer/
- Jay Alammar (2019), A Visual Guide to Using BERT for the First Time, http://jalammar.github.io/a-visual-guide-to-using-bert-for-the-first-time/
- NLP with Transformer, https://github.com/nlp-with-transformers/notebooks
- Min-Yuh Day (2022), Python 101, https://tinyurl.com/aintpupython101