# CS 410 Course Project – Final Report

## Source code and test set predictions

Code:
https://github.com/wfcwfcwfcwfc/CourseProject/blob/main/cs410-classification-contest.ipynb
Test set predictions:
https://github.com/wfcwfcwfcwfc/CourseProject/blob/main/cs410-classification-contest-result.txt

## Explain your model, and how you perform the training. Describe your experiments with other methods that you may have tried and any hyperparameter tuning.

The classifier uses BERT transformer deep learning framework at its core. BERT is a recent NLP framework based on Transformer and self-attention architecture. It serves the "encoder" in the transformer model and is widely adopted in text generation and text classification.

Training a BERT model generally has two stages: pre-training and fine-tuning. Pre-training aims at providing BERT a general understanding of a language. This step builds the embeddings and trains the parameters. Fine-tuning is optimizing BERT for certain specific tasks. Pre-training requires large language corpus and tremendous computing power. A common practice is to use existing pre-trained model and fine-tuning for the specific task. In this scenario, I use "bert-large-cased-whole-word-masking" from hugging-face as pre-trained model, then fine-tuned with the training data provided.

After fine-tuning, the model is capable to perform predictions. The performance with default parameters beats the baseline.

I also explored other pre-trained models like "distilbert-base-uncased", "roberta-base", "xlnet-base-cased". They all have smaller number of parameters compared to "bert-large-cased-whole-word-masking". The performance is good on training set but does not pass the baseline in test data.
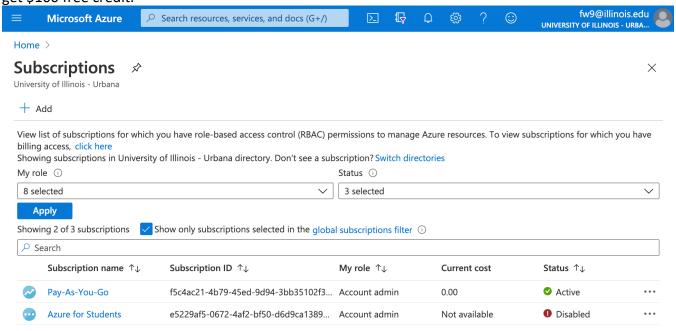
|  | distilbert-base-uncased | roberta-base | xlnet-base-cased |
|---|---|---|---|
| accuracy | 0.784800 | 0.753600 | 0.770400 |
| f1-score | 0.784064 | 0.751230 | 0.769815 |
| time(hrs) | 0.029202 | 0.038811 | 0.053048 |

On engineering side, the model was implemented with PyTorch and deployed on Microsoft Azure ML Studio. It provided convenient middleware for ML tasks for easy deployment and prototyping. All the fine-tuning was done on a single compute, GPU instance and running time is less than 10 minutes. Compute GPU instance pre-installed with CUDA 10.1. All code is contained in the notebook shown in the beginning of this document. I used the "NLP Best Practices" library as well as "NLP Utilities" library to build the classifier. These libraries also provided example templates which is referenced in this project.
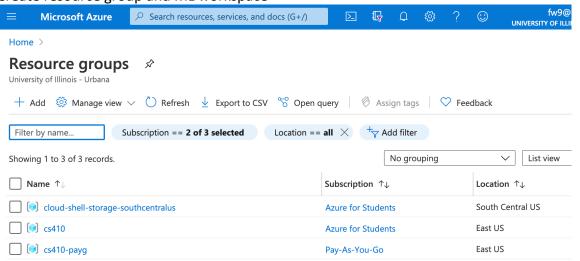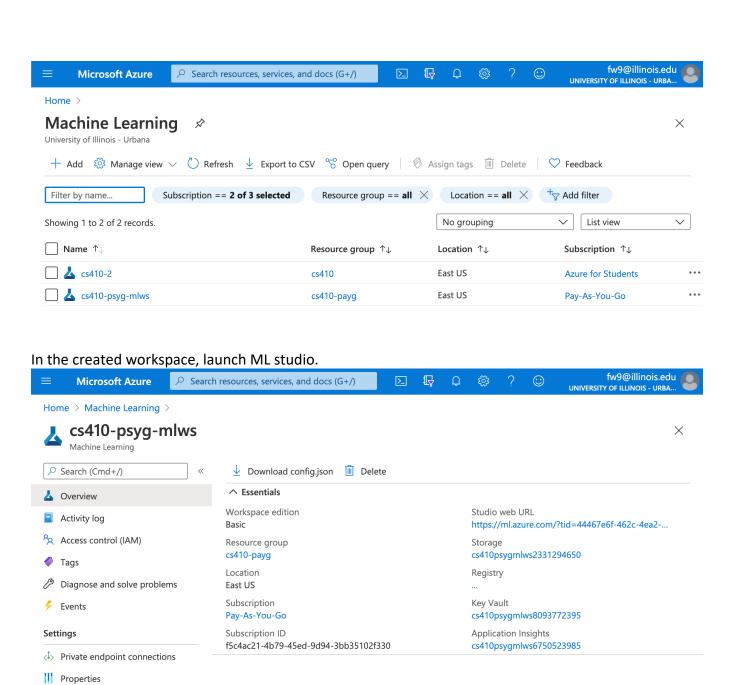
# Demo and Tutorial

## Environment Setup

Sign up for Microsoft Azure. Create a subscription that allows you to use GPU instances. Student email get $100 free credit.



## Create resource group and ML workspace

Home >

# Machine Learning 📌
University of Illinois - Urbana

✕

+ Add    ⚙ Manage view ∨    ↻ Refresh    ↓ Export to CSV    ⟲ Open query    |    🏷 Assign tags    🗑 Delete    |    ♡ Feedback

| Filter by name... | Subscription == **2 of 3 selected** | Resource group == **all** ✕ | Location == **all** ✕ | +🔽 Add filter |

Showing 1 to 2 of 2 records.

| No grouping ∨ | List view ∨ |

| ☐ Name ↑↓ | Resource group ↑↓ | Location ↑↓ | Subscription ↑↓ | |
|---|---|---|---|---|
| ☐ 🧪 cs410-2 | cs410 | East US | Azure for Students | ⋯ |
| ☐ 🧪 cs410-psyg-mlws | cs410-payg | East US | Pay-As-You-Go | ⋯ |

In the created workspace, launch ML studio.

Home > Machine Learning >

# 🧪 cs410-psyg-mlws
Machine Learning

✕

↓ Download config.json    🗑 Delete

∧ **Essentials**

| | |
|---|---|
| Workspace edition | Studio web URL |
| Basic | https://ml.azure.com/?tid=44467e6f-462c-4ea2-... |
| Resource group | Storage |
| cs410-payg | cs410psygmlws2331294650 |
| Location | Registry |
| East US | ... |
| Subscription | Key Vault |
| Pay-As-You-Go | cs410psygmlws8093772395 |
| Subscription ID | Application Insights |
| f5c4ac21-4b79-45ed-9d94-3bb35102f330 | cs410psygmlws6750523985 |

**Sidebar navigation:**

🧪 Overview
📋 Activity log
👥 Access control (IAM)
🏷 Tags
🔧 Diagnose and solve problems
⚡ Events

**Settings**
⟨⟩ Private endpoint connections
‖ Properties
🔒 Locks

**Monitoring**
🔔 Alerts
📊 Metrics

## Manage your machine learning lifecycle

Use the Azure Machine Learning studio to build, train,
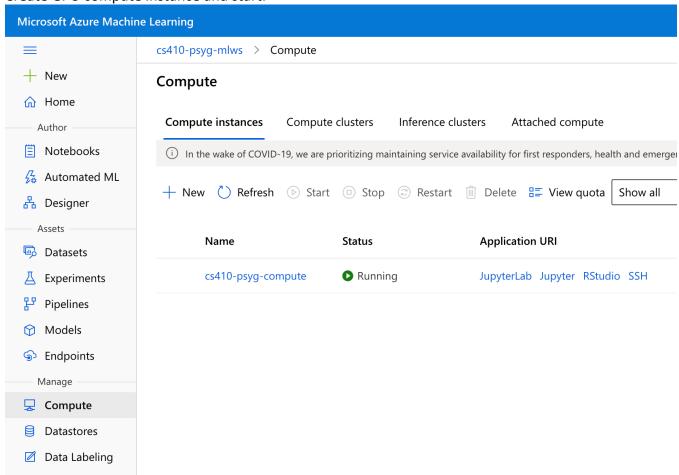evaluate, and deploy machine learning models. Learn more ↗

**Launch studio**

Create GPU compute instance and start.
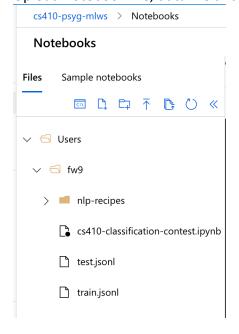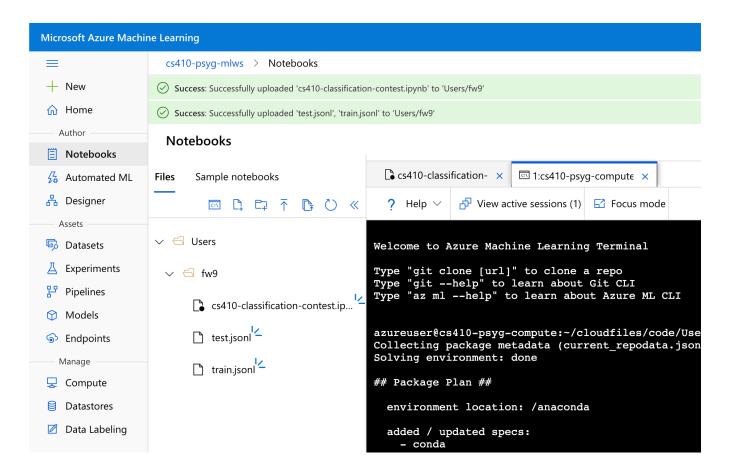
**Microsoft Azure Machine Learning**

| ☰ | cs410-psyg-mlws > Compute |
|---|---|

**＋ New**
**⌂ Home**

Author

📓 Notebooks
⚡ Automated ML
🔗 Designer

Assets

📊 Datasets
🔬 Experiments
🔀 Pipelines
📦 Models
☁ Endpoints

Manage

🖥 **Compute**
🗄 Datastores
✍ Data Labeling

## Compute

| Compute instances | Compute clusters | Inference clusters | Attached compute |
|---|---|---|---|

ⓘ In the wake of COVID-19, we are prioritizing maintaining service availability for first responders, health and emergen...

＋ New    ↻ Refresh    ▷ Start    ☐ Stop    ⟳ Restart    🗑 Delete    ☰ View quota    | Show all

| Name | Status | Application URI |
|---|---|---|
| cs410-psyg-compute | ▶ Running | JupyterLab  Jupyter  RStudio  SSH |

Upload notebook file, data file and clone the NLP Library.

cs410-psyg-mlws > Notebooks

## Notebooks

**Files**    Sample notebooks

⊞ 📄 📁 ⬆ 📄 ↻ ≪

∨ 📁 Users
  ∨ 📁 fw9
    > 📁 nlp-recipes
    📄 cs410-classification-contest.ipynb
    📄 test.jsonl
    📄 train.jsonl

## Update Conda:
```
conda update -n base -c defaults conda
```

## Generate conda env config:
```
cd nlp-recipes
python tools/generate_conda_file.py --gpu
```

Open nlp_gpu.yaml, update pytorch version from 1.4.0 to 1.5.1

## Create Anaconda env
```
conda env create -n nlp_gpu -f nlp_gpu.yaml
conda activate nlp_gpu
```

## Register this virtual env to notebook
```
python -m ipykernel install --user --name nlp_gpu --display-name "Python (nlp_gpu)"
```

```
The environment setup is complete at this point.
```

## Running the notebook
```
Open notebook and set kernel as 'nlp_gpu'.
Run the notebook, and the result shows up in "answer.txt"
```
Intermediate output can be seen in the notebook. Performance passed the baseline.

| 51 | fw9 | 6 | 0.6815589353612167 | 0.7966666666666666 | 0.7346311475409836 | 1 |