aws re: Invent

AIM 402-R

Deep Learning with PyTorch

Kris Skrinak

Sr. AI/ML Specialist PSA Amazon Web Services **Chaitanya Hazarey**

AI/ML Specialist PSA Amazon Web Services





Agenda

- State of the art in NLP using SageMaker
- SageMaker PyTorch Support
- Lab [1] Word Level Language Modeling Using PyTorch
- Train On Premise Deploy in the Cloud
- Lab [2] Toxic Comments Classification using Hugging Face Transformers and SageMaker
- Bonus Lab CAPTUM PyTorch Deep Learning Models Interpretability

Related sessions

AIM402-R1 [REPEAT] Deep learning with PyTorch

AIM412-R Deep learning applications using PyTorch, featuring Autodesk

AIM412-R1 Deep learning applications with PyTorch, featuring Freshworks

AIM407-R Amazon SageMaker and PyTorch: Tips & tricks

AIM407-R1 - [REPEAT] Amazon SageMaker and PyTorch: Tips & tricks

State of the art in NLP using SageMaker





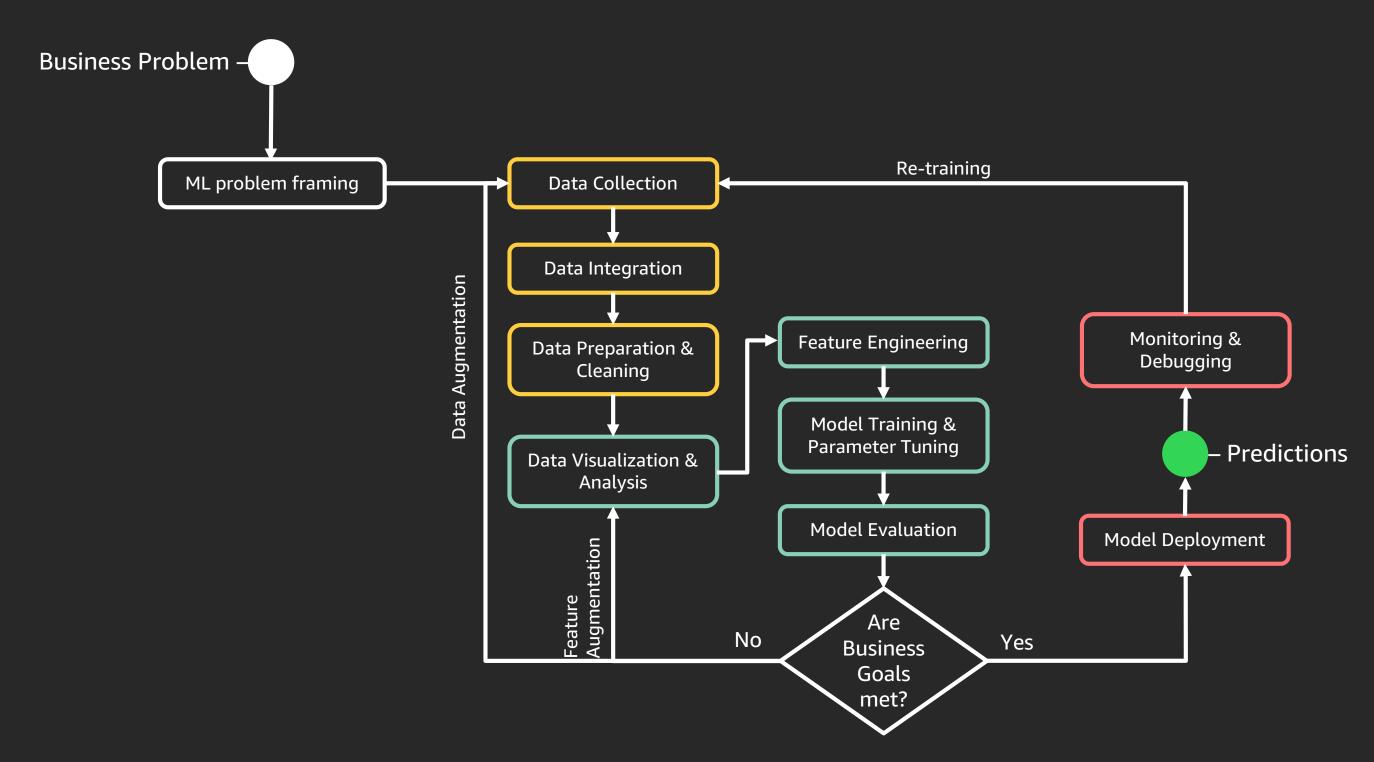
State of the art in NLP

- Essential components of NLP
 - Tokenize
 - Remove Stop Words
 - Lemmatize
 - Replace Rare Words
 - Embeddings
- Evolution of NLP from RNNs to Transformers
 - Seq2Seq
 - Attention
 - Transformer

Source:

- A Deep Dive into NLP with PyTorch PyData London 2019
- https://www.youtube.com/watch?v=4jROlXH9Nvc

The Machine Learning Process



THE AWS ML STACK

Broadest and deepest set of capabilities

AI Services

VISION			SPEECH		LANGUAGE		CHATBOTS FORECASTING		RECOMMENDATIONS	
REKOGNITION IMAGE	REKOGNITION VIDEO	TEXTRACT	POLLY	TRANSCRIBE	A 文文 TRANSLATE	COMPREHEND & COMPREHEND MEDICAL	LEX	FORECAST	PERSONALIZE	

ML Services

Amazon SageMaker Ground Trut	Notebooks	Algorithms + Marketplace	Reinforcement Learning	Training	Optimization	Deployment	Hosting	

ML Frameworks + Infrastructure

FRAMEWORKS	INTERFACES	INFRASTRUCTURE			
**TensorFlow mxnet	€ GLUON	Õ			
PYT <mark>Ö</mark> RCH	K Keras	EC2 P3 & P3DN	EC2 G4	E	



AMAZON SAGEMAKER

Bringing machine learning to all developers

Pre-built notebooks for common problems

Built-in, high performance algorithms

One-click training

One-click Optimization deployment

Fully managed with auto-scaling, health checks, automatic handling of node failures, and security checks



Collect and prepare training data



Choose and optimize your ML algorithm



Set up and manage environments for training



Train and tune model (trial and error)



Deploy model in production



Scale and manage the production environment

Launch SageMaker





https://dashboard.eventengine.run







Who are you?

- 1. By using Event Engine for the relevant event, you agree to the <u>AWS Event Terms and Conditions</u> and the <u>AWS Acceptable Use</u>
 <u>Policy</u>. You acknowledge and agree that are using an AWS-owned account that you can only access for the duration of the relevant event. If you find residual resources or materials in the AWS-owned account, you will make us aware and cease use of the account. AWS reserves the right to terminate the account and delete the contents at any time.
- 2. You will not: (a) process or run any operation on any data other than test data sets or lab-approved materials by AWS, and (b) copy, import, export or otherwise create derivate works of materials provided by AWS, including but not limited to, data sets.
- 3. AWS is under no obligation to enable the transmission of your materials through [AWS Event Engine] and may, in its discretion, edit, block, refuse to post, or remove your materials at any time.
- 4. Your use of the [event engine] will comply with these terms and all applicable laws, and your access to [AWS Event Engine] will immediately and automatically terminate if you do not comply with any of these terms or conditions.

Team Hash (e.g. abcdef123456)

This is the 12 digit hash that was given to you or your team.

Invalid Hash





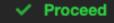


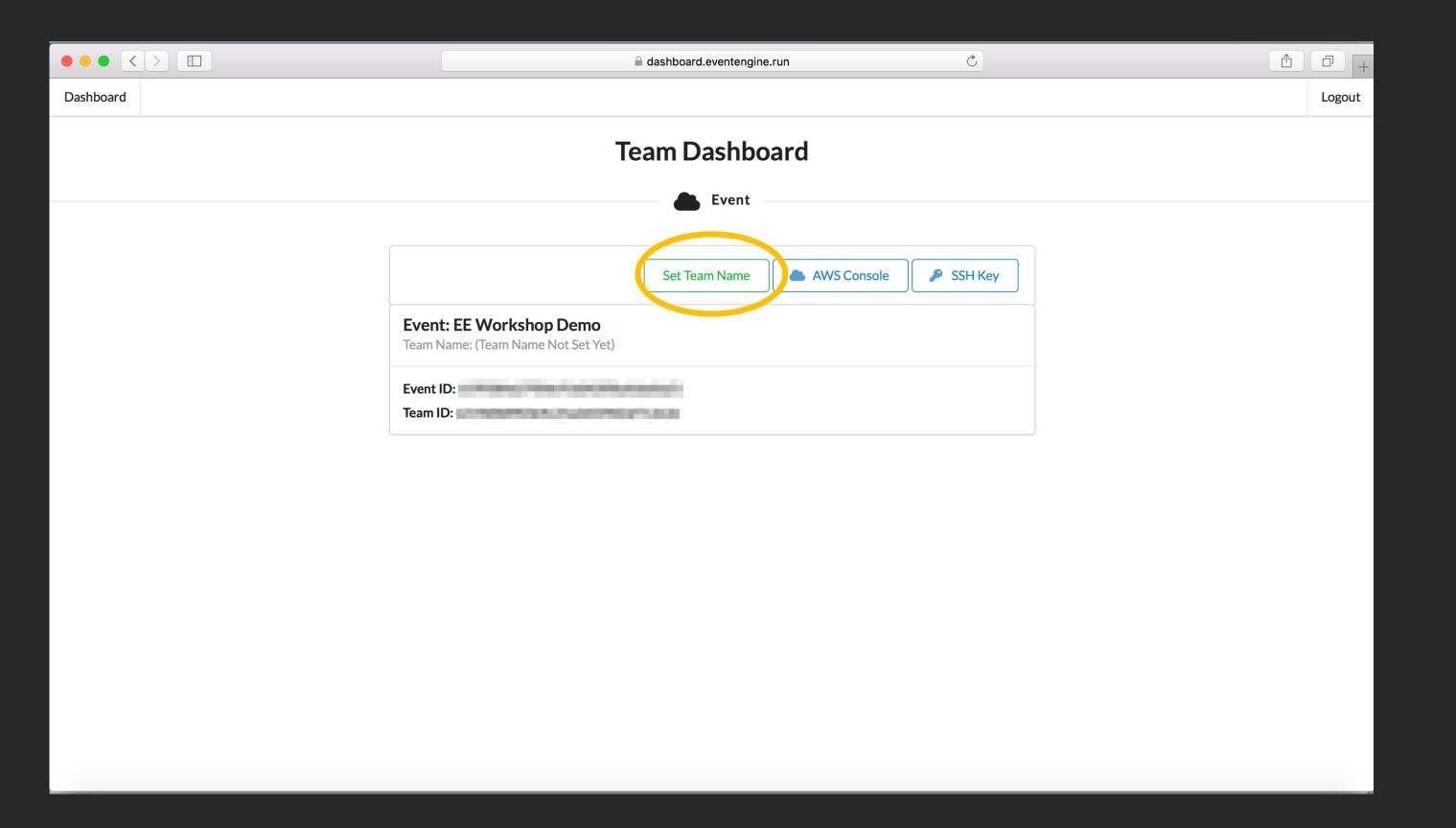
Who are you?

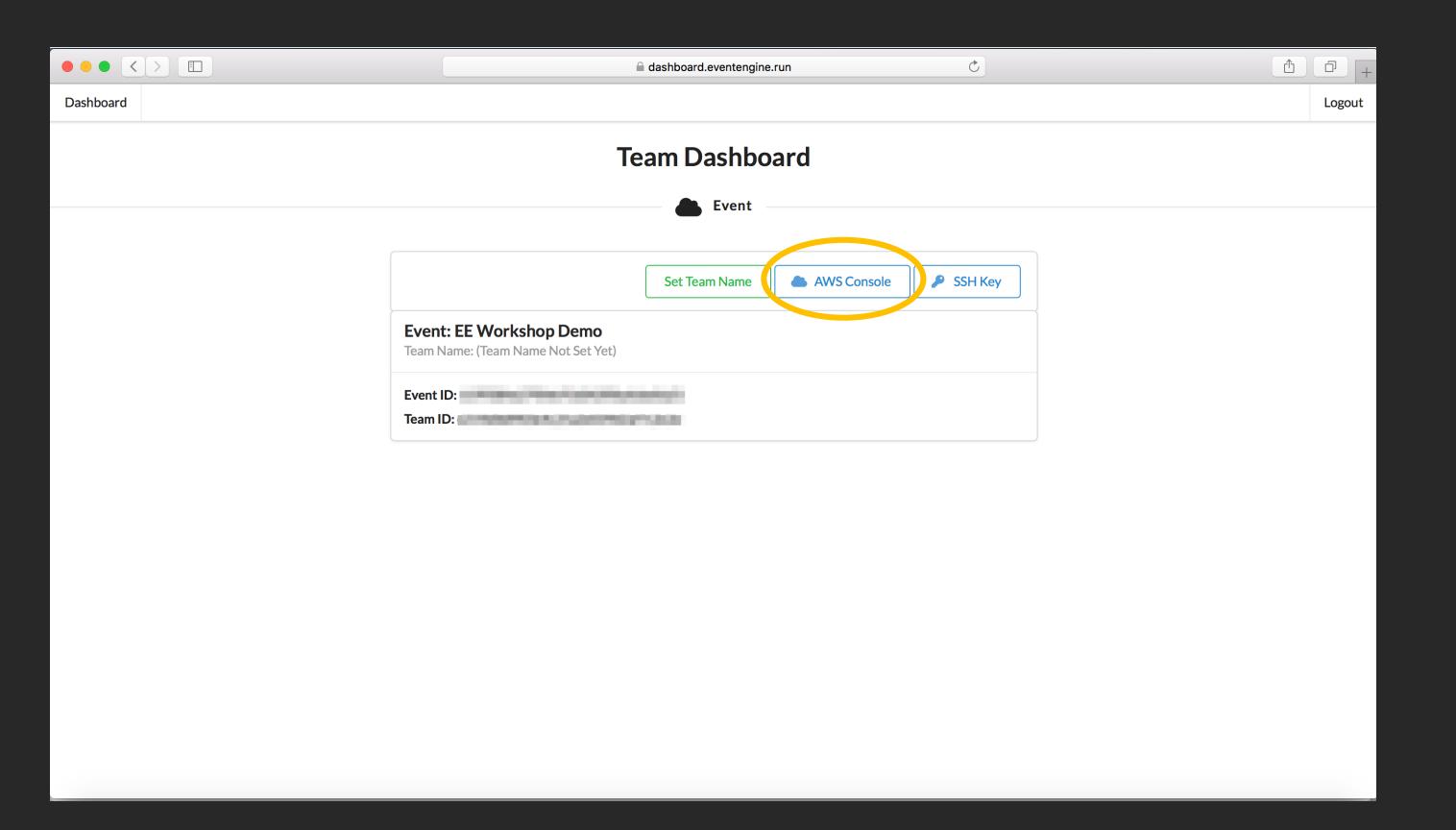
- 1. By using Event Engine for the relevant event, you agree to the <u>AWS Event Terms and Conditions</u> and the <u>AWS Acceptable Use Policy</u>. You acknowledge and agree that are using an AWS-owned account that you can only access for the duration of the relevant event. If you find residual resources or materials in the AWS-owned account, you will make us aware and cease use of the account. AWS reserves the right to terminate the account and delete the contents at any time.
- 2. You will not: (a) process or run any operation on any data other than test data sets or lab-approved materials by AWS, and (b) copy, import, export or otherwise create derivate works of materials provided by AWS, including but not limited to, data sets.
- 3. AWS is under no obligation to enable the transmission of your materials through [AWS Event Engine] and may, in its discretion, edit, block, refuse to post, or remove your materials at any time.
- 4. Your use of the [event engine] will comply with these terms and all applicable laws, and your access to [AWS Event Engine] will immediately and automatically terminate if you do not comply with any of these terms or conditions.



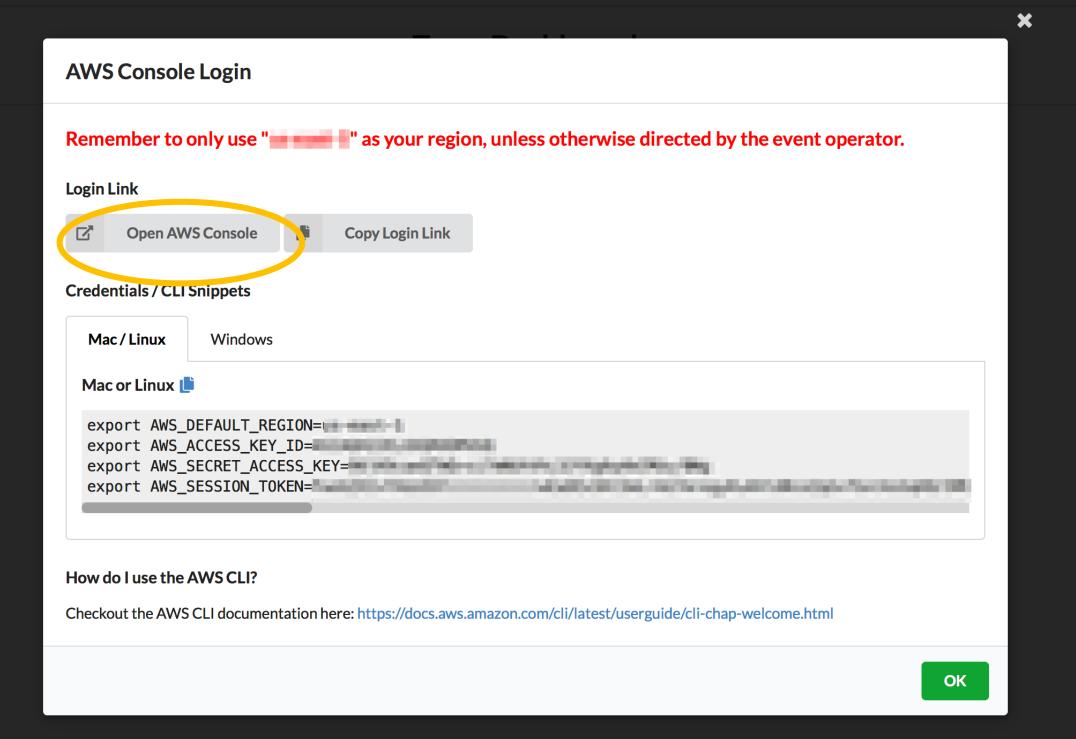
This is the 12 digit hash that was given to you or your team.

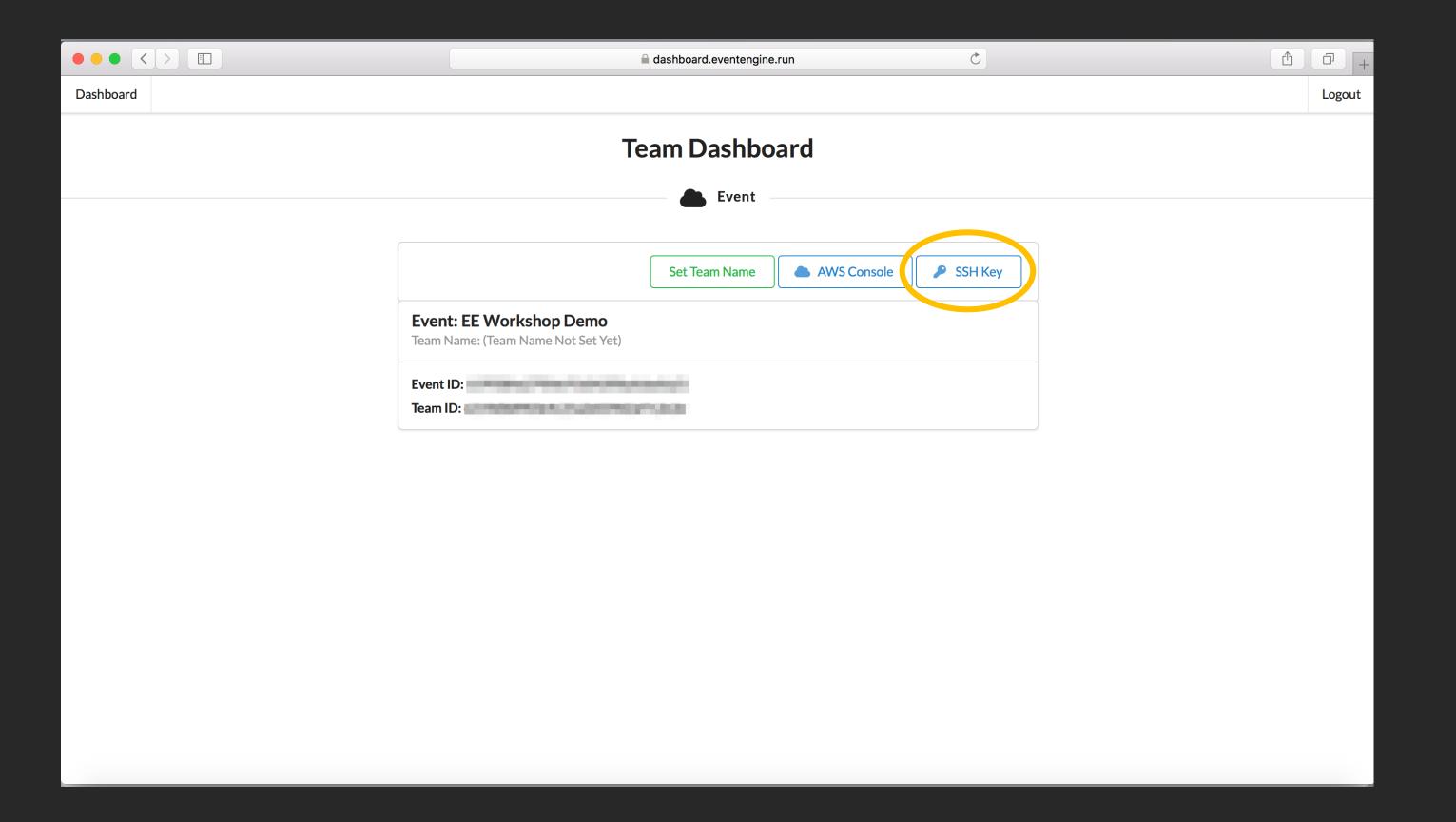






C



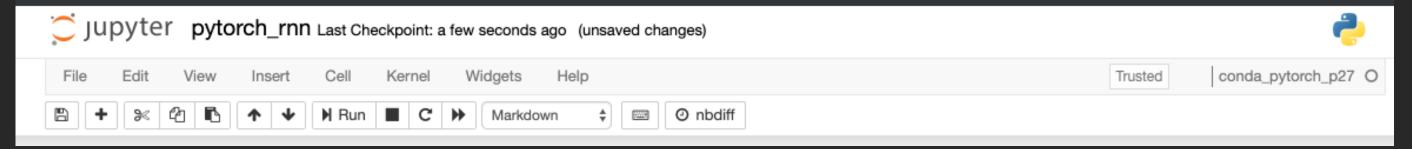


SageMaker PyTorch Support



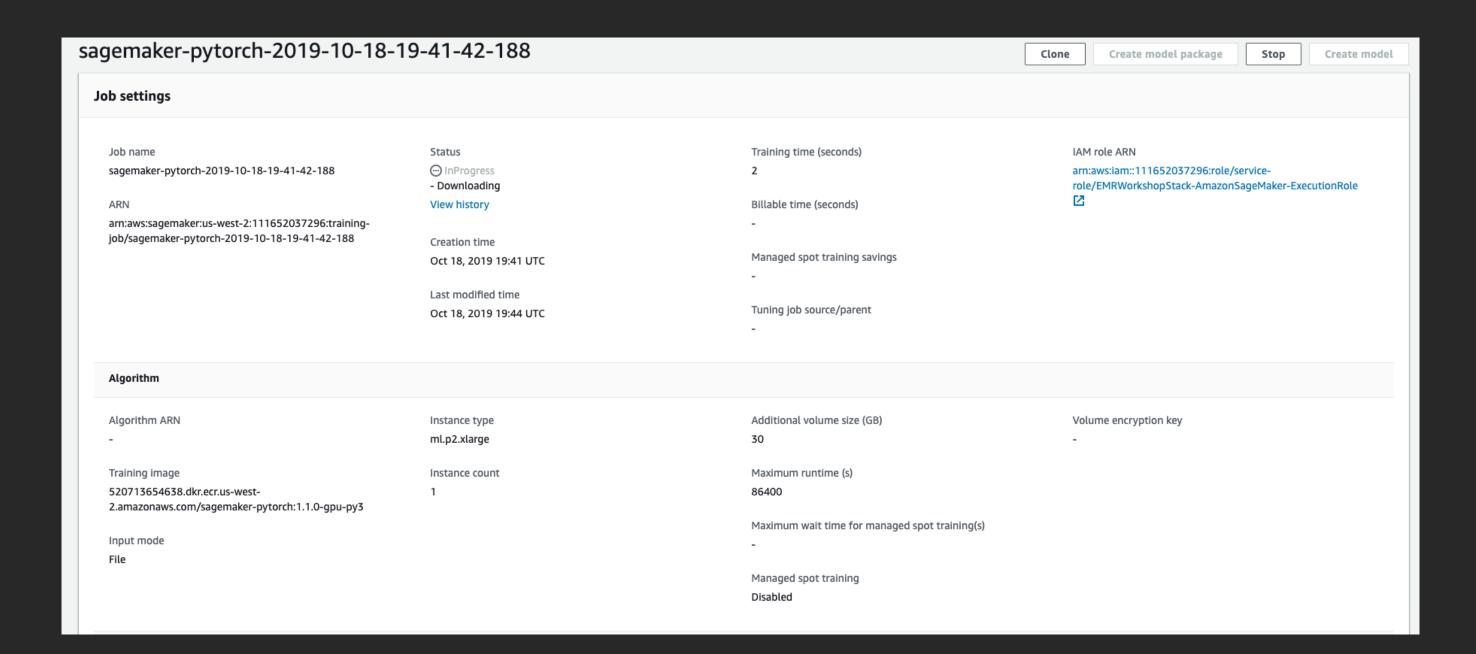


SageMaker PyTorch Support

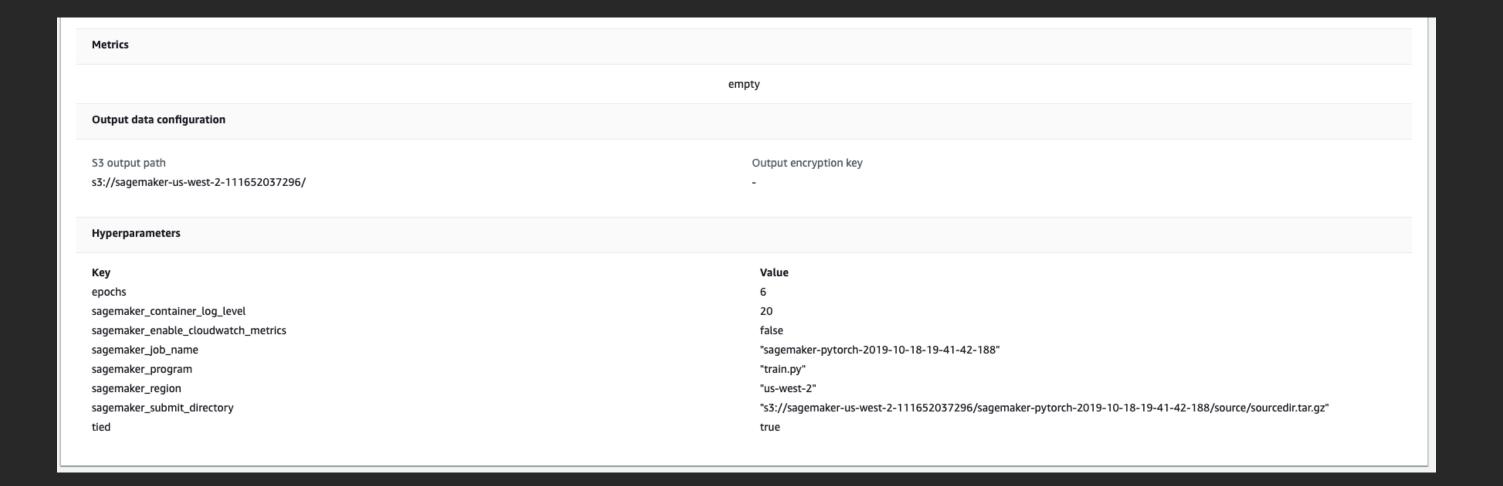


Estimator

Training Job



Hyperparameters



Training Script

```
    Jupyter train.py

✓ an hour ago

            View Language
                                                                                                                                Python
  1 # Based on github.com/pytorch/examples/blob/master/word language model
  2 import argparse
  3 import math
  4 import os
  5 from shutil import copy
  6 import time
    import torch
  8 import torch.nn as nn
 10 import data
 11 from rnn import RNNModel
parser = argparse.ArgumentParser(description='PyTorch Wikitext-2 RNN/LSTM Language Model')
 15 # Hyperparameters sent by the client are passed as command-line arguments to the script.
 parser.add_argument('--emsize', type=int, default=200,
                        help='size of word embeddings')
 18 parser.add argument('--nhid', type=int, default=200,
                        help='number of hidden units per layer')
 20 parser.add argument('--nlayers', type=int, default=2,
                        help='number of layers')
 22 parser.add argument('--lr', type=float, default=20,
                        help='initial learning rate')
 24 parser.add_argument('--clip', type=float, default=0.25,
 25
                        help='gradient clipping')
 26 parser.add_argument('--epochs', type=int, default=40,
                        help='upper epoch limit')
    parser.add_argument('--batch_size', type=int, default=20, metavar='N',
 29
                        help='batch size')
    parser.add_argument('--bptt', type=int, default=35,
 31
                        help='sequence length')
    parser.add argument('--dropout', type=float, default=0.2,
 33
                        help='dropout applied to layers (0 = no dropout)')
    parser.add argument('--tied', type=bool, default=False,
 35
                        help='tie the word embedding and softmax weights')
    parser.add_argument('--seed', type=int, default=1111,
                        help='random seed')
    parser.add_argument('--log-interval', type=int, default=200, metavar='N',
 39
                        help='report interval')
 40
 41 # Data and model checkpoints/otput directories from the container environment
    parser.add argument('--model-dir', type=str, default=os.environ['SM MODEL DIR'])
    parser.add_argument('--output-data-dir', type=str, default=os.environ['SM_OUTPUT_DATA_DIR'])
    parser.add argument('--data-dir', type=str, default=os.environ['SM CHANNEL TRAINING'])
 45
 46 args = parser.parse args()
```

Training Script

```
169
170 def train():
171
         # Turn on training mode which enables dropout.
172
         model.train()
173
         total loss = 0.
174
         start time = time.time()
175
         hidden = model.init hidden(args.batch size)
176
         for batch, i in enumerate(range(0, train data.size(0) - 1, args.bptt)):
177
             data, targets = get batch(train data, i)
             # Starting each batch, we detach the hidden state from how it was previously produced.
178
179
             # If we didn't, the model would try backpropagating all the way to start of the dataset.
180
             hidden = repackage hidden(hidden)
             model.zero grad()
181
182
             output, hidden = model(data, hidden)
183
             loss = criterion(output.view(-1, ntokens), targets)
184
             loss.backward()
185
186
             # `clip grad norm` helps prevent the exploding gradient problem in RNNs / LSTMs.
             torch.nn.utils.clip grad norm(model.parameters(), args.clip)
187
188
             for p in model.parameters():
                 p.data.add (-lr, p.grad.data)
189
190
191
             total loss += loss.item()
192
193
             if batch % args.log interval == 0 and batch > 0:
                 cur loss = total loss / args.log interval
194
195
                 elapsed = time.time() - start time
                 print('| epoch {:3d} | {:5d}/{:5d} batches | lr {:02.2f} | ms/batch {:5.2f} | '
196
197
                       'loss {:5.2f} | ppl {:8.2f}'.format(
198
                           epoch, batch, len(train data) // args.bptt, lr,
199
                           elapsed * 1000 / args.log interval, cur loss, math.exp(cur loss)))
200
                 total loss = 0
201
                 start time = time.time()
202
```

Model

```
97 # Build the model
    99
   print('Build the model')
100
101 ntokens = len(corpus.dictionary)
    rnn type = 'LSTM'
    model = RNNModel(rnn type, ntokens, args.emsize, args.nhid, args.nlayers, args.dropout, args.tied).to(device)
103
104
   criterion = nn.CrossEntropyLoss()
105
106
    # Save the data into model dir to be used with the model later
107
108 for file name in os.listdir(args.data dir):
       full_file_name = os.path.join(args.data_dir, file_name)
109
       if os.path.isfile(full file name):
110
111
          copy(full_file_name, args.model_dir)
112
    # Save arguments used to create model for restoring the model later
    with open(model_info path, 'wb') as f:
114
115
       model info = {
116
           'rnn type': rnn type,
117
           'ntoken': ntokens,
118
           'ninp': args.emsize,
119
           'nhid': args.nhid,
120
           'nlayers': args.nlayers,
121
           'dropout': args.dropout,
122
           'tie weights': args.tied
123
124
       torch.save(model info, f)
125
126
```

Save the model

```
print('Starting training.')
209
     for epoch in range(1, args.epochs+1):
210
         epoch start time = time.time()
211
         train()
212
         val loss = evaluate(val data)
213
         print('-' * 89)
214
         print('| end of epoch {:3d} | time: {:5.2f}s | valid loss {:5.2f} | '
215
               'valid ppl {:8.2f}'.format(epoch, (time.time() - epoch start time),
216
                                          val loss, math.exp(val loss)))
217
         print('-' * 89)
218
         # Save the model if the validation loss is the best we've seen so far.
219
         if not best state or val loss < best state['val loss']:</pre>
220
             best state = {
221
                 'epoch': epoch,
222
                 'lr': lr,
223
                 'val loss': val loss,
224
                 'val ppl': math.exp(val loss),
225
226
             print('Saving the best model: {}'.format(best state))
227
             with open(checkpoint path, 'wb') as f:
228
                 torch.save(model.state dict(), f)
229
             with open(checkpoint state path, 'w') as f:
230
                 f.write('epoch {:3d} | lr: {:5.2f} | valid loss {:5.2f} |
231
                         'valid ppl {:8.2f}'.format(epoch, lr, val loss, math.exp(val loss)))
232
         else:
233
             # Anneal the learning rate if no improvement has been seen in the validation dataset.
234
             lr /= 4.0
235
```

Model Deployment - PyTorchPredictor

- model_fn() load model
- input_fn() convert input payload into a PyTorch Tensor
- predict_fn() predict_fn() generates predictions from the model based on the return value of input_fn()
- output_fn() serializes the output from predict_fn() so that it can be returned by the SageMaker Endpoint

Hosting the model

```
In [ ]: predictor = model.deploy(initial_instance_count=1, instance_type='ml.m4.xlarge')
```

Lab 1 – Word-level Language modeling using PyTorch





Interpreting Results

Train on Premise – Deploy in Cloud using SageMaker





Train on premise – Deploy in the cloud

- Start the workflow from your laptop
- Use your favorite IDE
- Deploy to the cloud using SageMaker





One-click Deployment Fully managed with auto-scaling for 75% less

Lab 2 – Toxic Comments Classification using Hugging Face Transformers and SageMaker





Toxic Comment Classification Challenge

- Comments from Wikipedia human labeled
- 6 Levels 6 Labels Multiclass Classification
- Use Fastbert as a wrapper to use HuggingFace's Models
- Debug the lab notebook and run inference on your own text

Source:

- https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/data
- https://github.com/kaushaltrivedi/fast-bert/tree/master/sample_notebooks

Interpreting Results

Bonus Lab – CAPTUM – Model Interpretability





Thank you!

Kris Skrinak

skrinak@amazon.com

Chaitanya Hazarey

chazarey@amazon.com







Please complete the session survey in the mobile app.



