

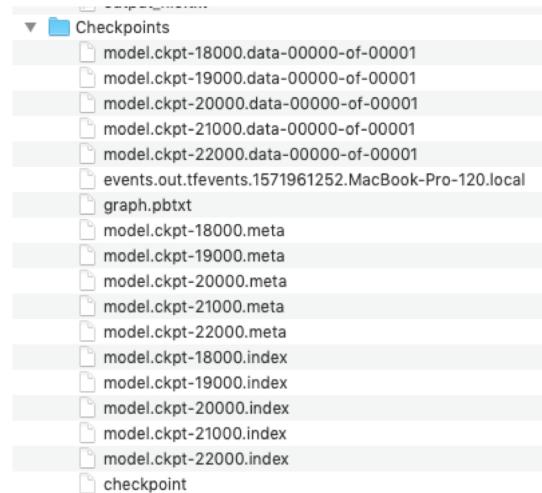
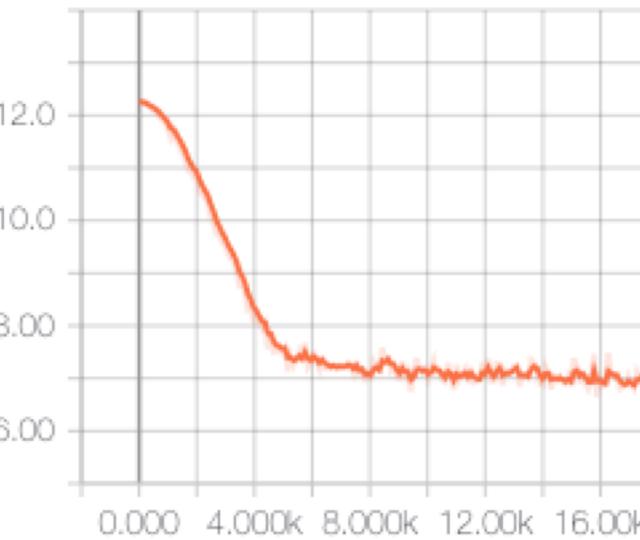
Experimenting with BERT Language Model

Al Gharakhanian

October 25, 2019

```
{  
    "attention_probs_dropout_prob": 0.1,  
    "directionality": "bidi",  
    "hidden_act": "gelu",  
    "hidden_dropout_prob": 0.1,  
    "hidden_size": 768,  
    "initializer_range": 0.02,  
    "intermediate_size": 3072,  
    "max_position_embeddings": 512,  
    "num_attention_heads": 12,  
    "num_hidden_layers": 12,  
    "pooler_fc_size": 768,  
    "pooler_num_attention_heads": 12,  
    "pooler_num_fc_layers": 3,  
    "pooler_size_per_head": 128,  
    "pooler_type": "first_token_transform",  
    "type_vocab_size": 2,  
    "vocab_size": 105879  
}
```

loss



e: Slimmer1401

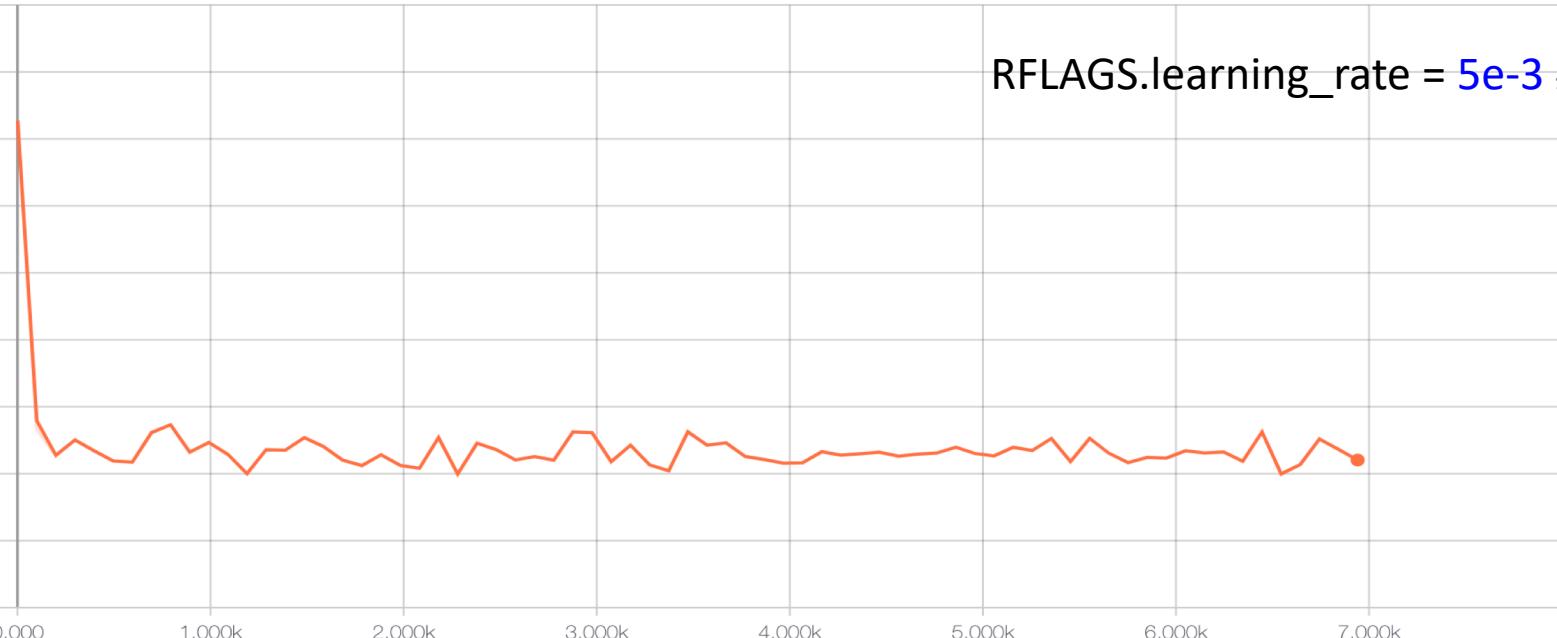
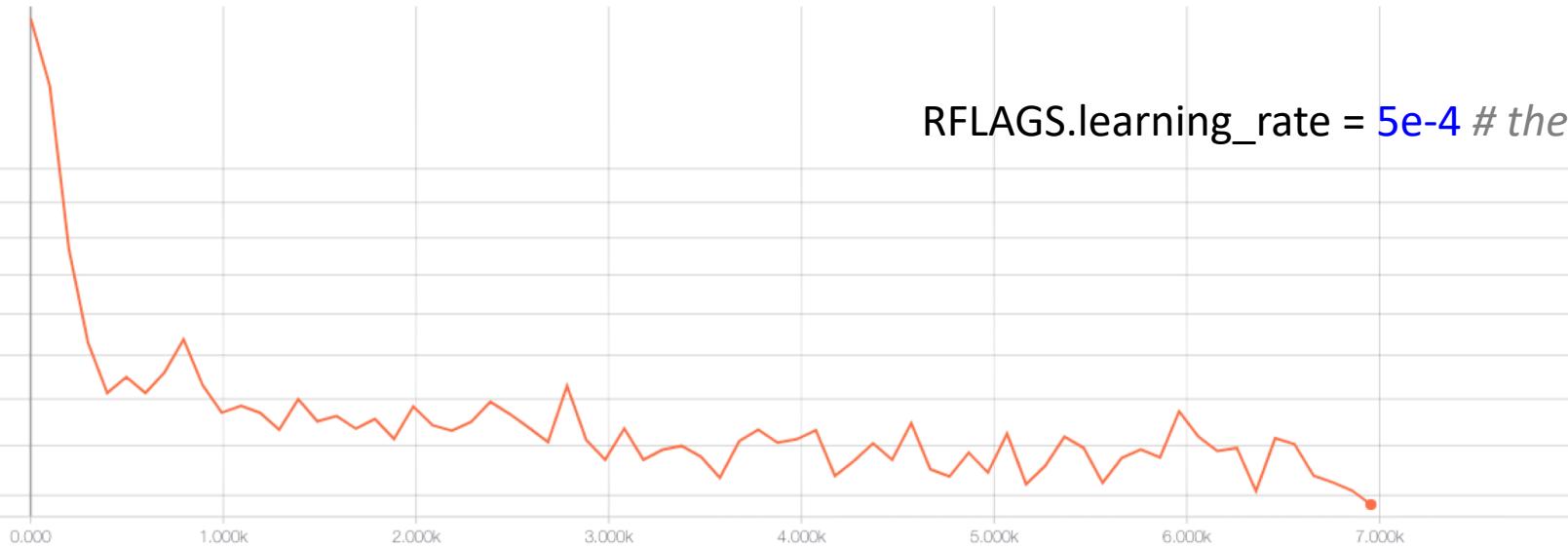
nt: Not sure about bias/variance since no CV

maybe learning rate too low (5000)

loss



```
0W:***** Eval results *****  
0W: global_step = 20000  
0W: loss = 7.0430512  
0W: masked_lm_accuracy = 0.076576576  
0W: masked_lm_loss = 6.545193  
0W: next_sentence_accuracy = 0.70125  
0W: next_sentence_loss = 0.51135576
```



```
RFLAGS.num_train_steps = 25000
```

```
10 # added this since the loss drops faster  
# the original was 5e-5
```



Multilingual vocab

English vocab

Run Name: Slimmer1401

Dataset Enhancements: None

Assessment:

RFLAGS.max_eval_steps = 1000

```
INFO:tensorflow:***** Eval results *****  
INFO:tensorflow: global_step = 25000  
INFO:tensorflow: loss = 4.831123  
INFO:tensorflow: masked_lm_accuracy = 0.30022892  
INFO:tensorflow: masked_lm_loss = 4.4728103  
INFO:tensorflow: next_sentence_accuracy = 0.8206875  
INFO:tensorflow: next_sentence_loss = 0.3701863
```

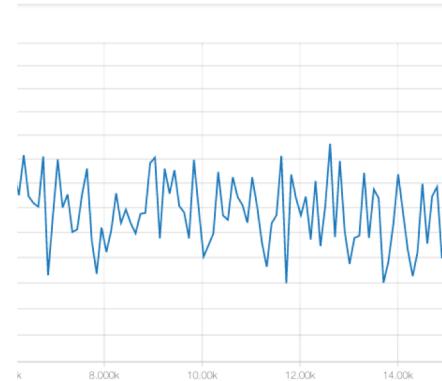
```
ropout_prob": 0.1,  
"bidi",  
"u",  
"ob": 0.1,  
  
": 0.02,  
": 384,  
ddings": 256,  
ds": 4,  
": 3,  
192,  
ion_heads": 4,  
ers": 2,  
ead": 128,  
rst_token_transform",  
2,  
79
```

Name:
immer1401max

dataset Enhancements: forced
ences to be max_seq_length
essment: the model has
rmed markedly worse in all
ts. Having a perfect next
nce accuracy means that there is

something wrong.

```
* Eval - re
lobal_step
ss = 5.89
sked_lm_a
sked_lm_l
ext_sente
ext_sente
```



```
"attention_probs_dropout_prob": 0.1,
"directionality": "bidi",
"hidden_act": "gelu",
"hidden_dropout_prob": 0.1,
"hidden_size": 96,
"initializer_range": 0.02,
"intermediate_size": 384,
"max_position_embeddings": 256,
"num_attention_heads": 4,
"num_hidden_layers": 3,
"pooler_fc_size": 192,
"pooler_num_attention_heads": 4,
"pooler_num_fc_layers": 2,
"pooler_size_per_head": 128,
"pooler_type": "first_token_trans",
"vocab_size": 2105879
"vocab_size": 105879
```

RFLAGS.max_eval_steps = 1000

: ad
e o
..

RFLAGS.max_eval_steps = 1000

```
INFO:tensorflow:***** Eval results *****
INFO:tensorflow: global_step = 25000
INFO:tensorflow: loss = 4.827713
INFO:tensorflow: masked_lm_accuracy = 0.30101895
INFO:tensorflow: masked_lm_loss = 4.4638457
INFO:tensorflow: next_sentence_accuracy = 0.817375
INFO:tensorflow: next_sentence_loss = 0.3772189
```

RFLAGS.max_eval_steps = 200

```
INFO:tensorflow:***** Eval results *****
INFO:tensorflow: global_step = 25000
INFO:tensorflow: loss = 4.8605146
INFO:tensorflow: masked_lm_accuracy = 0.29909173
INFO:tensorflow: masked_lm_loss = 4.481621
INFO:tensorflow: next_sentence_accuracy = 0.818125
INFO:tensorflow: next_sentence_loss = 0.38994846
```

```
INFO:tensorflow:***** Eval results *****
INFO:tensorflow: global_step = 25000
INFO:tensorflow: loss = 4.831123
INFO:tensorflow: masked_lm_accuracy = 0.30022892
INFO:tensorflow: masked_lm_loss = 4.4728103
INFO:tensorflow: next_sentence_accuracy = 0.8206875
INFO:tensorflow: next_sentence_loss = 0.3701863
```

RFLAGS.max_eval_steps = 100

```
INFO:tensorflow:***** Eval results *****
INFO:tensorflow: global_step = 25000
INFO:tensorflow: loss = 4.9306307
INFO:tensorflow: masked_lm_accuracy = 0.2951858
INFO:tensorflow: masked_lm_loss = 4.539931
INFO:tensorflow: next_sentence_accuracy = 0.805
INFO:tensorflow: next_sentence_loss = 0.41647342
```

```
:tensorflow: global_step = 25000
:tensorflow: loss = 5.04019
:tensorflow: masked_lm_accuracy = 0.28756413
:tensorflow: masked_lm_loss = 4.6136303
:tensorflow: next_sentence_accuracy = 0.8219375
:tensorflow: next_sentence_loss = 0.4329269
```

2000

```
:tensorflow: global_step = 25000
:tensorflow: loss = 5.054097
:tensorflow: masked_lm_accuracy = 0.28611273
:tensorflow: masked_lm_loss = 4.625028
:tensorflow: next_sentence_accuracy = 0.82040626
:tensorflow: next_sentence_loss = 0.43505368
```

4000

```
tensorflow:***** Eval results *****
tensorflow: global_step = 25000
tensorflow: loss = 5.0459895
tensorflow: masked_lm_accuracy = 0.28678566
tensorflow: masked_lm_loss = 4.6231885
tensorflow: next_sentence_accuracy = 0.821975
tensorflow: next_sentence_loss = 0.4283578
```

25000

loss



- Run Name:

Slimmer1401_delimited

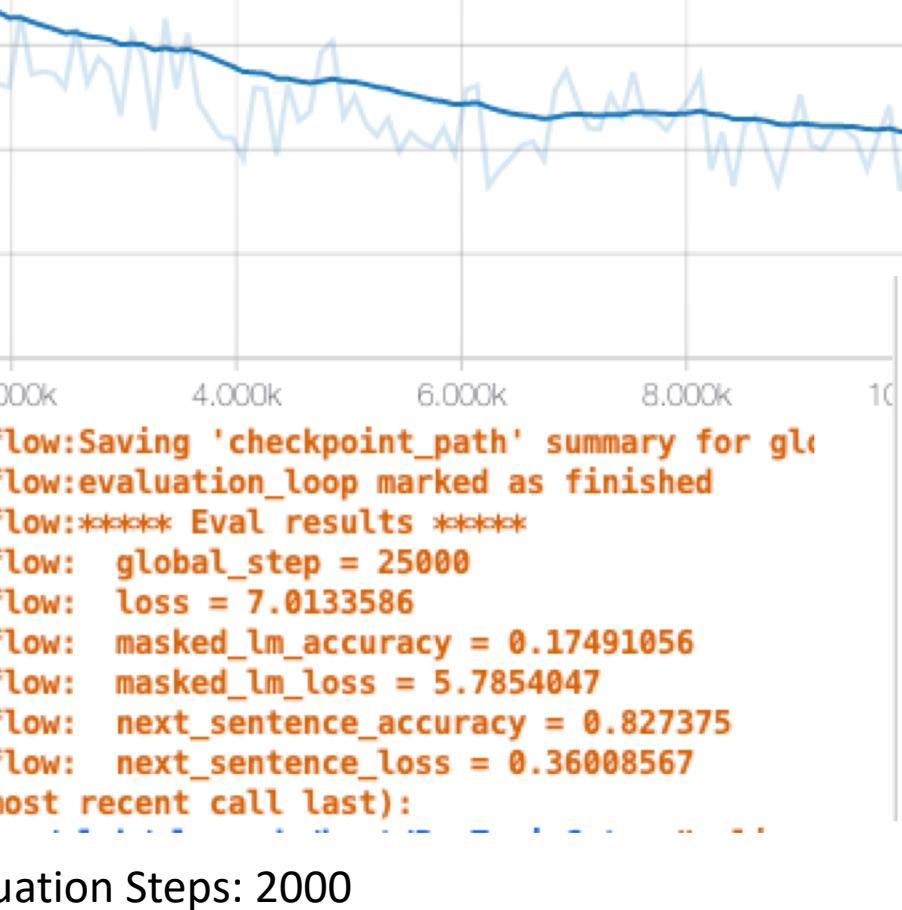
- Dataset Enhancements: **None**

- Description: In this version the examples are aligned based on '' sentence ending

- Assessment: The model did worse than the baseline

RFLAGS max_eval_steps = **1000**
in every metrics category

- Observation: the loss curve seems to be much chopier



```
INFO:tensorflow: loss after : [6.09838724]
INFO:tensorflow:Evaluation [25000/25000]
INFO:tensorflow:Finished evaluation at 2019-12-01-06:16:54
INFO:tensorflow:Saving dict for global step 25000: global_step=25000, loss=6.09838724, masked_lm_accuracy=0.17482635, masked_lm_loss=5.7844195, next_sentence_accuracy=0.824615, next_sentence_loss=0.36632505
INFO:tensorflow:***** Eval results *****
INFO:tensorflow: global_step = 25000
INFO:tensorflow: loss = 7.001802
INFO:tensorflow: masked_lm_accuracy = 0.17482635
INFO:tensorflow: masked_lm_loss = 5.7844195
INFO:tensorflow: next_sentence_accuracy = 0.824615
INFO:tensorflow: next_sentence_loss = 0.36632505
all done
```

Evaluation Steps: 25000

Impact of Dummy Token Insertion in BERT Pretraining

Description	Parameters/Environment	Results
<p>The intent of this experiment is to determine whether inserting dummy tokens in the sentences of the training data, used to pretrain BERT language model; improves the overall performance of the pretraining. The rationale is to make the context recognition more robust which can lead to enhanced performance in downstream tasks. The dummy tokens are selected randomly from the prior two sentences. Single characters and fragments are excluded in the selection. Total number of intrusions in two modes are roughly equivalent. Note: in order to reduce the training time, the size of the model is reduced markedly. Multitask learning has been employed to accommodate the insertion loss implementation.</p>	<p>Environment Python 3.6 / Tensorflow 1.11 Dataset Project Gutenberg Word Embedding Byte Pair Encoding Learning Rate 5.00E-04 Training Steps 70000 Evaluation Steps 20000 Warm Up Steps 2000 Num. of Inserted Tokens 5 Masked Token Probability 15% without insertion / 10% with insertion</p>	<p>The insertion of dummy tokens in the training examples does not impact the long term total loss during latter stages of the training. The most notable observation of this experiment is that insertion of dummy tokens significantly accelerates the drop of the total loss during the early phase of the training. It turns out that the model is extremely effective in discerning between a valid token vs. an inserted token which explains why the impact of insertion loss dissipates quickly. Next step is to examine the impact of the insertion on the performance of downstream tasks such as reading comprehension (SQuAD).</p>

