CS 505: Introduction to Natural Language Processing

> Wayne Snyder Boston University

Lecture 19 – Sampling Strategies for Generative Models; Transformer Architecture Continued; Pretrained Models; BERT



Lecture Plan

- Sampling in Generative Models (esp. HW 05)
- o Transformer Architecture Continued
 - Encoder Review
 - Decoder Architecture
 - Transfer Learning
 - BERT and friends....

Some Strategies for Sampling in Generative Models

In Generative Models, "Sampling" refers to the the choice of the next token using the probability distribution output by the network.

Let's look at a <u>really simple</u> example:

Suppose the generative model is creating sentences from the vocabulary

V = ['<s>', '</s>', 'NLP', 'early', 'fun', 'green', 'hard', 'interesting', 'is']

We have generated three tokens so far:

['<s>', 'NLP', 'is']

and the softmax output (probability distribution) for the next token is:

<u>Token</u>	Probability
<s></s>	0.0
	0.01
NLP	0.03
early	0.12
fun	0.12
green	0.0
hard	0.2
interesting	0.51
is	0.01

Some Strategies for Sampling in Generative Models

There are five common strategies for sampling.

['<s>', 'NLP', 'is', ???]

1. Sampling using the unmodified distribution:

1 print('\nsample: ',sample_choice(outcomes,distribution,option=0)) 0.0 <s> </s> 0.01 NLP 0.03 early 0.12 fun 0.12 0.0 green hard 0.2 interesting 0.51 0.01 is sample: fun

from numpy.random import choice

choice(a=outcomes,p=distribution)

Code for sample_choice(...) is on web site.

2. Greedy Sampling: Just choose the single most probable:

1	print('\nsample:	',sample_choice(outcomes,distribution,option=1))
<s></s>		0.0	
<td>></td> <td>0.0</td> <td></td>	>	0.0	
NLP		0.0	
ear	ly	0.0	
fun		0.0	
gree	en	0.0	
ĥaro	d	0.0	
inte	erestin	q 1.0	
is		0.0	
samp	ole: i	nteresting	

Some Refinements for Sampling in Generative Models



['<s>', 'NLP', 'is', ???]

1 print('	<pre>\nsample: ',sample_choice(outcomes,distribution,option=2,temperature = 0.01))</pre>			
<s></s>	7.095474162284459e-23			
	1.928749847963851e-22			
NLP	1.4251640827408859e-21			
early	1.1548224173015387e-17			
fun	1.1548224173015387e-17			
green	7.095474162284459e-23			
hard	3.4424771084698575e-14		Temperature	Distribution
interesting	0.999999999999955		Temperature	Distribution
is	1.928749847963851e-22	Cold: Exaggerate		
		Cold. Exaggerate	0	One Hot
sample: in	teresting	differences	U U	enerier
		1		
: 1 print('\ns	<pre>ample: ',sample_choice(outcomes,distribution,option=2,temperature = 0.1))</pre>			
			10	Lisual Softmax
<s></s>	0.005446284988308871		1.0	Usual Solullax
	0.0060190757806309345			
NLP	0.0073517157600377	•		
early	0.018082302955730302	Hat Minimiza		
fun	0.018082302955730302	Hot: Minimize		
green	0.005446284988308871	differences	8	Uniform
hard	0.040242905309378116			
interesting	0.893310051481244			
is	0.0060190757806309345			

sample: interesting

1 print('\n	<pre>sample: ',sample_choice(outcomes,distribution,option=2,temperature = 0.5))</pre>	1 print('\n	<pre>sample: ',sample_choice(outcomes,distribution,option=2,temperature = 10.0)</pre>
<s> </s> NLP early fun green hard interesting is	0.08396476946896886 0.08566097032727966 0.08915686084440715 0.10674014184435929 0.10674014184435929 0.08396476946896886 0.12526071682556345 0.23285065904881377 0.08566097032727966	<s> </s> NLP early fun green hard interesting is	0.10986989484510404 0.10997981969321279 0.11019999943895172 0.11119627595335972 0.10986989484510404 0.11208941394957388 0.11561860562812123 0.10997981969321279
sample:		sample: gree	n

Some Refinements for Sampling in Generative Models

4. Top-K Sampling: From top K most probable, choose proportionally:

1 print('\n	<pre>sample: ',sample_choice(</pre>	outcomes,distribution,option=3(,K=4))		
interesting hard fun early sample: fun	0.5368421052631579 0.2105263157894737 0.12631578947368421 0.12631578947368421	Normalize to a probability distribution (sum = 1.0)	interesting hard fun early NLP is green <s></s>	0.51 0.2 0.12 0.03 0.01 0.01 0.0

5. Top-P: From top choices with sum of probability p, choose proportionally:

1 print('\n	<pre>sample: ',sample_choice(ou</pre>	utcomes,distribution,option=4,p=0.75))			<u>Sum</u>
interesting hard	0.6144578313253012 0.24096385542168675		interesting hard fun	0.51 0.2 0.12	0.51 0.71 0.83
fun sample: inte	0.14457831325301204 resting	Normalize to a probability distribution (sum = 1.0)	early NLP is	0.12 0.03 0.01	0.95 0.98 0.99
			 green <s></s>	0.01 0.0 0.0	1.0

Some Refinements for Sampling in Generative Models

Code for sample choice (...) is on web site.

```
1 # Options
2
3 #
        0 = No modification (default)
4 # 1 = Greedy: return one-hot of maximum prob, equivalent to softmax with 0 temperature
5 # 2 = Softmax with temperature (temperature of 1.0 is ordinary softmax)
6 # 3 = Top-K: Choose from top K most probable options (no other modification)
7 # 4 = Nucleus (Top-p): p is a cutoff, choose from the top choices whose cumulative prob just
8 #
                exceeds p (so in general, cumulative prob is slightly more than p)
9
10
11 import math
12 import numpy as np
13 from numpy.random import choice
14
15 def sample_choice(outcomes,distribution,option=3,temperature=0.3,K=5,p=0.25):
16
17
       if type(outcomes) != np.ndarray:
18
           outcomes = np.array(outcomes)
19
       if type(distribution) != np.ndarray:
20
           distribution = np.array(distribution)
21
22
23
       if option == 0:
24
           return choice(a=outcomes,p=distribution/sum(distribution))
25
       elif option == 1:
26
           oneHot = np.array([0.0]*len(distribution))
           oneHot[np.argmax(distribution)] = 1.0
27
28
           return outcomes[np.argmax(distribution)]
29
       elif option == 2:
           sum_exp = sum(math.exp(x/temperature) for x in distribution)
30
31
           lst = [math.exp(x/temperature)/sum_exp for x in distribution]
32
           return choice(a=outcomes,p=[math.exp(x/temperature)/sum exp for x in distribution])
33
       elif option == 3:
34
           indices = np.argsort(distribution)[-K:]
35
           return choice(a=outcomes[indices],p=distribution[indices]/sum(distribution[indices]))
36
       elif option == 4:
37
           indices = np.argsort(distribution)
38
           total = 0
39
           for C in range(len(distribution)-1,-1,-1):
40
               total += distribution[indices[C]]
41
               if total > p:
42
                   break
43
           indices = indices[C:]
44
           return choice(a=outcomes[indices],p=distribution[indices]/sum(distribution[indices]))
45
```

Transformer: Multi-Head Attention

The center of the design is the Attention mechanism, here Multi-Head Attention....



Transformer: Multi-Head Attention

The basic mechanism at work here is self-attention, where the input is processed to determine the dependencies between different words.

Self-attention is implemented by a series of linear transformations, scaled, and then softmaxed to produce the probability distribution which tells us how much each word depends on other words in the same sequence.





Scaled Dot-Product Attention

Transformer: Multi-Head Attention

Multi-Head Attention applies this self-attention mechanism as 8 (or more) self-attention Heads:



This allows the encoder to try to understand 8 different kinds of dependencies among the words in the input.



The original design applied 8 self-attention heads to sequences of length 512.

512 words

Transformer: Encoder

Then, this encoder layer is stacked 6, 8, or more layers deep!





Transformer: Decoder

The Decoder stage has essentially the same features as the Encoder, except for a stage called Masked Multi-Head Attention.

Apparently, the classic diagram (from the original paper) is a bit misleading, and there are additional connections between the encoder and decoder:



Transformer: Decoder

The Decoder stage has essentially the same features as the Encoder, except for Multi-Headed Attention is **masked** to simulate left to right generation of output; Each stage of Decoder takes the context vector (output of encoder) as an additional input:



Transfer Learning: Leveraging pre-trained models

A hugely important development in ML in the last decade is the development of transfer learning systems, consisting of two stages:

Pre-Training: A (large) model is trained on a (large) corpus to produce generic outputs;

Downstream: A NN is added on top of the pre-trained model to build models for a different task (called a "downstream task").



Transfer Learning: Leveraging pre-trained models

There are two flavors of transfer learning:

Feature-based: Use outputs of pre-trained model as "features" input to downstream model; further training only occurs in downstream model:

Fine-Tuning: Outputs of pre-trained model are inputs to downstream model, but BOTH models are trained; typically, changes to pre-trained model are minimal.



Example: ELMO

Examples: GPT, BERT

Transfer Learning: Leveraging pre-trained models

This was first used in image processing (but is now used everywhere):



AlexNet structure of transfer learning from ImageNet database. The parameters are transferred in all layers from Conv1 to Conv5 except FC.

Transfer Learning for NLP: GPT, ELMO, and BERT

Pre-trained models are typically very large and trained for a long time on huge amounts of data (so they have a lot of knowledge you can leverage for your specific task.

The most successful systems in NLP are

- GPT (Generative Pre-trained Transformer)
- ELMO (Embeddings from Language Models)
- BERT (Bidirectional Encoder Representations from Transformers)

These all

- Are based on transformers, but use ONLY the encoder phase;
- Produce contextual embeddings;
- GPT is autoregressive (outputs feed back into inputs)





Transfer Learning for NLP: GPT, ELMO, and BERT

The most significant difference in the models is how they process the input sequence:

GPT uses a simplified "causal" transformer model which only uses left-context:



Figure 11.1 A causal, backward looking, transformer model like Chapter 9. Each output is computed independently of the others using only information seen earlier in the context.

Transfer Learning for NLP: GPT, ELMO, and BERT

The most significant difference in the models is how they process the input sequence:

ELMO and BERT use a bidirectional transformer architecture:



Figure 11.2 Information flow in a bidirectional self-attention model. In processing each element of the sequence, the model attends to all inputs, both before and after the current one.

BERT

1 - Semi-supervised training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp

2 - Supervised training on a specific task with a labeled dataset.



BERT

- Bert has two implementations:
- Bert base:
- 12 layers; 12 attention heads per laye
- 768 hidden units; 110 M parameters
- Bert large:
- 16 layers; 16 attention heads per laye
- 1024 hidden units; 340 M parameters

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	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERTLARGE	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

12

2

1





BERT added two significant ideas for training which allowed it to achieve SOTA performance on significant tasks:

Training using a Masked Language Model (MLM); and

Focus on a "next sentence prediction" task with two sentences as input.

Note: The first versions of BERT used the 800 M word BooksCorpus and a 2.5 B word English Wikipedia corpus.



The Masked Language Model is based on an educational theory/testing paradigm known as the Cloze Task, where students learn a language by filling in blanks in a story or piece of text:

	Word vault	>
Once upon a time there was an pig with	X old	
hree pigs, and because there was	× little	
not enough to them, she sent them	X feed	
nut to their fortunes	🗙 seek	
	× first	
he little pig went off and met a man	🗙 straw	
who had a bundle of . He asked the	X straw	
man "Please give me that to build a	X house	
	🗙 straw	

Masked Language Modeling uses unannotated text from a large corpus. 15% of the words in the corpus are selected for the training phase: of these,

80% of replaced with the token [MASK]10% are replaced with randomly-selected tokens10% are left unchanged.

The model is trained to predict the missing tokens



A variation of MLM uses spans (subsequences of the input sequence); all of the words in the span are replaced as before:



Figure 11.6 Span-based language model training. In this example, a span of length 3 is selected for training and all of the words in the span are masked. The figure illustrates the loss computed for word *thanks*; the loss for the entire span is based on the loss for all three of the words in the span.

The Next Sentence Prediction task is to input TWO sentences starting with the token [CLS] and separated by the token [SEP]. The training set has 50% sentences that are next to each other in the corpus, and 50% random sentences.



Using BERT

Bert can be used for sentence classification if a single sentence is input:



Figure 11.8 Sequence classification with a bidirectional transformer encoder. The output vector for the [CLS] token serves as input to a simple classifier.

BERT can classify the relationship between two sentences:

• Neutral

- a: Jon walked back to the town to the smithy.
- b: Jon traveled back to his hometown.
- Contradicts
 - a: Tourist Information offices can be very helpful.
 - b: Tourist Information offices are never of any help.
- Entails
 - a: I'm confused.
 - b: Not all of it is very clear to me.



BERT can generate the most likely sentence to follow a given sentence:



Use Beam Search to find most likely sentence to follow.

Using BERT

Bert can be used for sequence labelling if all of the outputs are used:



Figure 11.9 Sequence labeling for part-of-speech tagging with a bidirectional transformer encoder. The output vector for each input token is passed to a simple k-way classifier.

BERT Punches Above Its Weight!

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	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERTBASE	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERTLARGE	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

