Practical Machine Learning Dr. Ashish Tendulkar Department of Computer Science and Engineering Indian Institute of Technology, Madras

Lecture - 30 Recurrent Neural Networks (Part 2)

So, let us try to understand how the loss is computed in case of simple RNN.

(Refer Slide Time: 00:27)



We use the input, perform non-linear activation over linear combination to obtain the output for each position.

So, we get the error at every position. So, for example, the loss at t^{th} position is a function of the actual value and the predicted value.

$$loss^{} = f(y^{}, \hat{y}^{}) = -y^{}log(\hat{y}^{}) - (1 - y^{})log(1 - \hat{y}^{})$$

In case of classification we can use cross entropy loss. So this is the loss that we incur, this is a cross entropy loss that we incur at each position in the sequence. So, the total loss of the total loss due to the prediction can be obtained simply by summing across the sequence length.

$$L(y, \hat{y}) = \sum_{t=1}^{T_x} loss^{}(y^{}, \hat{t}^{})$$

So, we obtained this loss at each position and in order to calculate the gradient of the loss with respect to the input. We use an algorithm called back propagation. So, in case of RNN; it is called as back propagation through time or BPTT.

So, each of these outputs are through this link dependent on each of the input variables and the outputs of all the previous states. So, we use back propagation through time to calculate the weights that minimize this particular loss. We know that simple RNNs suffer from vanishing gradient problem and we also know that there is a neat trick that was used to make changes in the architecture so that vanishing gradient problem can be addressed.

So, let us look at the modified RNN architecture that takes care of vanishing gradient problem; it is called as LSTM or Long Short Term Memory.

(Refer Slide Time: 06:51)



So, beginning with our simple RNN. We had this architecture for simple RNN where we gave a feature vector corresponding to each position as input. We obtain output for each layer and in addition to the input feature vector, we have a recurrent connection which is the output from the previous node that is also used as an additional input for calculating the output.

In order to solve the vanishing gradient problem faced by the simple RNNs because they cannot really remember information from some previous node. We add a provision to carry information from previous node and use that as an additional input while calculating the output.

So, now in case of simple RNN, we just had the input vector of the particular position and the recurrent connection. In case of LSTM, we have an additional input which is carry which is some information from previous states. Now we have to figure out how to decide which information to carry on the carry line. Let us try to understand how to decide that particular information. We will be using the same strategy as simple RNN for deciding what to carry in the next state. Let us understand that through a pseudo code.

(Refer Slide Time: 10:03)



So, the output will be computed based on the dot-products of the output from the previous state or the recurrent connection along with its weight; the input vector and its rate, the carry information in its weight and the bias term. We add all these dot-products and biases and apply non-linear activation on that to get the output at t^{th} position. In order to decide what to carry further, we define three different non-linear activations on some dot-products.

$$output_t = activation(dot(output_t, U_0) + dot(input_t, W_0) + dot(c_t, V_0) + b_0)$$

So, we define the first output i_t which sums up the output from the previous state or from previous timestamp and its feature vector. Then this is the linear combination or this is a dot-product between the feature vector of that position and its weight and the bias term; f_t is calculated by applying non-linear activation on the dot-product between the output of the previous state and its feature vector and dot-product between the feature vector and its weight vector and the bias, k t is again calculated in the same manner.

$$i_t = activation(dot(state_t, U_i) + dot(input_t, W_i) + b_i)$$

$$f_t = activation(dot(state_t, U_f) + dot(input_t, W_f) + b_f)$$

$$k_t = activation(dot(state_t, U_k) + dot(input_t, W_k) + b_k)$$

So, all these three numbers are calculated based on activation applied to dot-product between the recurrent connection or the recurrent information and its weights; adding that into a dot-product between the feature vector and its weight, adding that to the bias term. And we decide a next carry is by multiplying i_t with k_t and c_t with f_t. c_t is the carry coming in and this is the formula for deciding the carry for the next stage.

$$c_t + 1 = i_t * k_t + c_t * f_t$$

So, let us try to demonstrate this in a picture. We will insert a small block to calculate carry. So, here we calculate the new carry and that carry is passed on the carry line and we also pass the state information coming from this box to the next box. So, these are architectural details. Most of the time we find such kind of architectures by doing search in the architecture space.



We will take this LSTM model and try to use this in practice to solve some of the problems related to text and time series. So, this is the LSTM model; the only change that we make is we use a carry information and we perform additional computations to decide what should be carried to the next state. We learn the weights of feature vectors used in all this calculation through the training of these LSTM models.

You can see that because of these different operations LSTM have far more number of parameters as compared to simple RNN. But LSTM models are quite powerful and they are showing state of the art results on lots of sequence learning tasks. Let us experience LSTMs in action by going through some of the practical examples.

(Refer Slide Time: 17:15)



So, let us use RNN models to obtain sentiment of movie reviews. In this example we will use movie reviews from IMDB and we will use RNN classifier to predict the sentiment of each of the movie reviews. The output here is binary; the review can either be positive or negative. We begin by importing the necessary libraries and downloading the movie review data set using tensorflow datasets.

(Refer Slide Time: 18:15)

<pre>************************************</pre>	()	+ Text & Copy to Drive	V RAM Cisk	· / Editing
<pre>Preserview 1.x selected. Import metplotlib and create helper function to plot graphs: [1] import metplotlib and create helper function to plot graphs: [2] import metplotlib.psplot as plt</pre>	٣	* Temperfuc_verion only exist in Colab. Temperfuc_verion 1.x except Exception: pass import temperfuc_datasets as tfds import temperfuc_as tff		NPTEL
<pre>Import satpletils and create a helper function to plot graphs: [] import satpletils.oyplet as plt def glot_graphs(hitory, string); plt_blot(hitory, string); plt_blot(hitory, interfact); plt_blot(hitory, interfact);</pre>	Θ	TensorFlow 2.x selected.		
<pre>[1] import metplotlik.opylet as pit def ging_probe(history, string); pit_pit(history, history("state"); pit_pit(history, history("state"); pit_history("state"); pit_history("state"); pit_history("state"); pit_history("state"); distate(interpine); pit_history("state"); distate("state");</pre>	Imp	ort matplotlib and create a helper function to plot graphs:		
<pre>efe gier_grephe(hitsoy, string): pit.gle(hitsoy, hitsoy)(string) pit.gle(hitsoy, hitsoy)(string) pit.gle(hitsoy)(string) pit.gle(hitsoy)(string)(string)(string) pit.gle(hitsoy)(string)(</pre>	[2]	import matplotlib.pyplot as plt		
Setup input pipeline The MADD large mover review dataset is a binary classification dataset—all the reviews have either a positive or negative sentiment. Download the dataset using <u>IFOS</u> The[dataset comes with an inbuil subword tokenizer. [1] dataset, info = tfds.load('indu',reviews('bubberd&B', with_infortrue,		<pre>def plot_graphs(history, string): pls.plot(history, sitrong) pls.plot(history, sitrong'("string]) pls.plot(history, sitrong'("st_"+string]) pls.plot(history, sitrong'("st_"+string]) pls.show() </pre>		
The BMDB large movie review dataset is a binary classification dataset—all the reviews have either a positive or negative sentiment. Download the dataset using <u>ITDS</u> The[dataset comes with an inbull subword takenizer. [3] dataset, info = tfids.load('imbg.reviews/ubwordsBL', with_infortTrue,	+ Se	tup input pipeline		
Download the dataset using IEDS The[dataset comes with an inbull subword takenizer. [3] dataset, info = tfds.load('im0b_reviews/subwordsRi', with_inforTrue, at_updrvisedTrue) train_dataset, test_dataset = dataset['true'].dataset('test'] Downloading and preparing dataset indb_reviews (RB.3 RIB) to /root/tensorflow_datasets/imdb_reviews/subwordsRi/0.1.0 DI Completed.	The	IMDB large movie review dataset is a binary classification dataset-all the reviews have eith	er a positive or negative sentiment.	
 [3] dataset, info * tfds:load('isst,reviews/vbbordsR', kith_inforfrue, as_uperisedFrue) train_dataset, test_dataset = dataset['this'], dataset['test'] @ Downloading and preparing dataset infb_reviews (80.23 RHB) to /root/tensorflow_datasets/imfb_reviews/subwordsR/0.1.0 Di Compleme Di Compleme Di Compleme 	Dov	moad the dataset using TFDS. The dataset comes with an inbuilt subword tokenizer.		
Downloading and preparing dataset indb_reviews (80.23 HiB) to /root/tensorflow_datasets/indb_reviews/subwordsBk/0.1.0 D Completed. Townloading and preparing dataset indb_reviews (80.23 HiB) to /root/tensorflow_datasets/indb_reviews/subwordsBk/0.1.0 D Completed. Townloading and preparing dataset indb_reviews (80.23 HiB)	[3]	<pre>dataset, info = tfds.load('imb_reviews/submodslk', with_infonTrue,</pre>		
Di Completed	0	Downloading and preparing dataset imdb_reviews (80.23 MiB) to /root/tens	orflow_datasets/imdb_reviews/subwords8k/0.1.0	
PL Cine ROUTE OF THE ROUTE OF		Di Completed 1//100% 1/1 [00:05<00:00, 5.59s/ url]		
bilote outprove evide (evidence volue, reve Alibia)		DI Size 801/100% 80/80 (00:05<00:00, 14.45 MB/	4	

This data set comes with an inbuilt sub word tokenizer.

(Refer Slide Time: 18:24)

3	+ Text & Copy to Drive	PAM Disk FAM Fditing
1	(i) Could not be transformed and mass of executed as-ast paese report that to the succession of a second process of the succession of t	Name was reading the bog, set the versionary to be (or
As th	is is a subwords tokenizer, it can be passed any string and the tokenizer will tokenize it.	
[4]	tokenizer = info.features['text'].encoder	
[5]	<pre>print ('Vocabulary size: {}'.foreat(tokenizer.vocab_size))</pre>	
θ	Vocabulary size: 8185	
[6]	<pre>sample_string = 'TensorFlow is cool.'</pre>	
	<pre>tokenized_string = tokenizer.encode(sample_string) print ('Tokenized_string is {}'.format(tokenized_string))</pre>	
	<pre>original_string = tokemizer.decode(tokemized_string) print ('The original string: {)'.format(original_string))</pre>	
	assert original_string == sample_string	
0	Tokenized string is [6307, 2327, 4043, 4265, 9, 2724, 7975] The original string: Tensorflow is cool.	
The t	okenizer encodes the string by breaking it into subwords if the word is not in its dictionary.	
[7]	<pre>for ts in tokenized string: print ('()> [)'.format(ts, tokenizer.decode([ts])))</pre>	
θ	6307> Ten	
	4043> F1	
	4265> Ou	
	2724 ····> cool	
	7975> .	

(Refer Slide Time: 18:27)

3	+ Text d Copy to Drive	V RAM
2	Tokenized string is [6387, 2327, 4043, 4265, 9, 2724, 7975] The original string: Temsorflow is cool.	NPTEL
The	tokenizer encodes the string by breaking it into subwords if the word is not in its dictionary.	
[7]	<pre>for ts in tokenized string: print ('{}> ()'.format(ts, tokenizer.decode([ts])))</pre>	
Θ	6387 ****> Ten 2327 ****> 567 4363 ****> 51 4365 ****> 6u 9 ***>> 1s 2724 ***> cool 2925 ****> .	
[8]	RUFFER_SIZE = 10000 BATCH_SIZE = 64	
[9]	train_detaxet = train_detaxet.bxdffd=(NUFRE_SIII) train_detaxet = train_detaxet.pxddxd_batch(NATO(SIII, train_detaxet.output_bhages) test_detaxet = test_detaxet.pkddxd_batch(NATO(SIII, test_detaxet.output_bhages)	
- Cre	rate the model	
Build	I a tf.kers. Sequential model and start with an embedding layer. An embedding layer stores one vector per word. When called, it converts the Indices to sequences of vectors. These vectors are trainable. After training (on enough data), words with similar meanings often have similar ve	e sequences of ctors.
This	index-lookup is much more efficient than the equivalent operation of passing a one-hot encoded vector through a tf.keras.layers.Dense layer	t.
A rec the r	surrent neural network (RNN) processes sequence input by iterating through the elements. RNNs pass the outputs from one timestep to their input ent.	ut-and then to
The	ed loans lawses Bidianstians) wannes can dee he wad with an 2000 laws. This economises the input forward and hadiousele through the	DNM Issuer and

After downloading the data set, we will shuffle the data set and we use padded_batch function to obtain the training data set where each sequence has a fixed length. We also used pattern_batch method to obtain the test data set.

(Refer Slide Time: 19:01)

	Text d Copy to Drive	Disk I	/ Editing
())	rrain_dataset = train_dataset.shuffle(BUFFELSIZE) rain_dataset = train_dataset.padde_batch(DATCH_SIZE, train_dataset.output_shapes)		N
	test_dataset = test_dataset.padded_batch(BATCM_SIZE, test_dataset.output_shapes)		
Crea	te the model		
Build word	Ef. karas. Sequential model and start with an embedding layer. An embedding layer stores one vector per word. When called, it converts the soc does to sequences of vectors. These vectors are trainable. After training (on enough data), words with similar meanings often have similar vector.	juences of s,	
This in	dex-lookup is much more efficient than the equivalent operation of passing a one-hot encoded vector through a tf.keras.layers.Dense layer.		
A recu the ne	rent neural network (RNN) processes sequence input by iterating through the elements. RNNs pass the outputs from one timestep to their input- t.	nd then to	
The to	.keras.layer.8idirectional wrapper can also be used with an RNN layer. This propagates the input forward and backwards through the RNN noatenates the output. This helps the RNN to learn long range dependencies.	layer and	
[10]	<pre>model = tf.kersi.Sequential(ff.kersi.Sequential(fi.kersi.Japers.Editectional(tf.kersi.Lapers.ISTM(64)), tf.kersi.Japers.Dense(64, activations'right'), tf.kersi.Japers.Dense(64, activations'signeds')))</pre>		
Comp	e the Keras model to configure the training process:		
[11]	<pre>model.compile(loss*binary_crossentropy',</pre>		
0	node].surkary()	↑ ↓ 00	01

Let us come to the model creation part; we are going to use RNN models over here. Here we use a bi-directional model on top of LSTM; the bi-directional wrapper propagates the

input forward and backward through the RNN layers and then concatenates the output. This helps RNN to learn long range dependencies.

So, we take the input which is text and pass it through the embedding layer which gets us a vector for each word. We embed each word into 64 length vector. We pass the output of embedding to bi-directional layer. So, the output of bi-directional layer is passed to a dense layer with 64 units and we use radio as an activation function. The output of dense layer is passed through the output layer which is again a dense layer with a single unit. As we have binary classification problem here and we use sigmoid as an activation function here. Let us define the model, compile it and look at model summary.

(Refer Slide Time: 20:55)

1	+ Text 🙆 Copy to Drive			Cisk Cisk Cisk
1	<pre>model = tf.keras.Sequential([tf.keras.layers.Embedding] tf.keras.layers.Bidirectic tf.keras.layers.Dense(64, tf.keras.layers.Dense(1, 4])</pre>	<pre>(tokenizer.vocab_size, nal(tf.keras.layers.L activation='relu'), activation='signoid')</pre>	64), TR((()),	
Comp	alle the Keras model to configure the	e training process:		
[11]	model.compile(loss='binary_cro optimizer='adam' metrics=['accura	<pre>ssentropy', (y'])</pre>		
0	model.summary()			<u>↑↓∞¢∎</u>
θ	Model: "sequential"			
	Layer (type)	Output Shape	Paran #	
	bidirectional (Bidirectional	(None, 128)	66848	
	dense (Dense)	(None, 64)	8256	
	dense_1 (Dense) Total params: 598,209 Trainable params: 598,209 Non-trainable params: 0	(None, 1)	65	
Trai	in the model			

So, you can see that embedding outputs 64 numbers and the bi-directional LSTM output 64 numbers for each direction: forward 64, backward 64. So, concatenation of that results into output containing 128 numbers. The dense layer output 64 numbers as we are using 64 units here and the final layer outputs a single number.

(Refer Slide Time: 21:46)



We start with text we have embedding layer followed by bi-directional LSTM followed by our dense layer and one more dense layer which is an output layer. We have let us say an input sequence of length *t*. So, let us say these are all LSTM units. So, here we are passing the recurrent connection in this particular direction from left to right; we actually get 64 numbers from the LSTM.

We are essentially passing that output to the next level and we are collecting the outputs only at the last layer. In the other pass, we start with last word and do LSTM calculations and pass the recurrence from right to left; in a sense because we started with the last word and we are going up to the first word and then both these outputs kind of concatenated and we get 128 units from this.

(Refer Slide Time: 25:07)



We start with LSTM model, then we take an input. So, we have LSTM model that outputs 64 number for input at each position. So, this is the forward pass because the output of t^{th} position or i^{th} position is being used as a recurrent connection in express i^{th} position. For example, the output of $x^{<1>}$ is being used as a recurrent connection for $x^{<2>}$.

So, this is called the forward pass in bi-directional one; we also define some kind of a backward pass, where you pass the result of i^{th} node to $i - 1^{th}$ node as a recurrent connection. For example, for the second position; the output of the second position will be passed back to the first position and will be used as a recurrent input. And we concatenate the outputs of let us at first position; the forward pass and the backward pass.

So, since here we are outputting 64 values each; the concatenation outputs 128 values. So, we perform this concatenation at each position for example, these two will get concatenated or these two will get concatenated and each one of them will output 128 numbers.

(Refer Slide Time: 29:19)

ode	+ Text			V RAM Disk - V Editing
[10]	<pre>model = tf.keras.Sequenti tf.keras.layers.Embed tf.keras.layers.Eldir tf.keras.layers.Dense tf.keras.layers.Dense])</pre>	al([ding(tokenizer.vocab_size, (ectional(tf.keras.layers.LS1 (64, activation='relu'), (1, activation='sigmoid')	4), M(64)),	
Corr	pile the Keras model to configu	are the training process:		
[11]	model.compile(loss='binar optimizer=' metrics=['a	y_crossentropy', adam', ccuracy'])		
0	model.summary()			↑↓∞¢∎
θ	Model: "sequential"			
	Layer (type)	Output Shape	Param #	
	embedding (Embedding)	(None, None, 64)	523840	
	bidirectional (Bidirecti	ional (None, 128)	66048	36
	dense (Dense)	(None, ∯4)	8256	
	dense_1 (Dense)	(None, 1)	65	A
	Total params: 598,209 Trainable params: 598,20	19		1 7 -2

And is 128 numbers are passed to dense layer which has got 64 units.

(Refer Slide Time: 29:31)

ode	+ Text			V RAM Disk 🔤 V Editing
[12]	bidirectional (Bidirectional (8	None, 128)	66848	
θ	dense (Dense) (1	None, 64)	8256	
	dense_1 (Dense) () Total params: 598,209 Trainable params: 598,209 Non-trainable params: 0	None, 1)	65	
Tra	in the model			
				↑↓∞¢∎
9	history = model.fit(train_dataset validation_data	t, epochs=2, ata=test_dataset)		↑↓∞¢≣
0 []	history = model.fit(train_dataset validation_dataset test_loss, test_acc = model.evalu	t, epochs=2, ata=test_dataset) uate(test_dataset)		^ ↓ ∞ ‡ ∎
0	history = model.fit(train_dataset validation_dat test_loss, test_acc = model.evalu print('Test_loss: ()'.format(test print('Test Accuracy: ()'.format)	rt, epochsal, ata=test_dataset) uate(test_dataset) t_loss)) (test_acc))		<u>↑↓∞¢</u>
C I I	history = model.fit(train_dataset validation_dat test_loss, test_acc = model.eval print("Test Loss: ()'.format(test print("Test Accuracy: ()'.format) above model does not mask the paddim	rt, epochs=2, iata=test_dataset) uate(test_dataset) t_loss)) (test_acc)) g applied to the sequ	ences. This can lead to skewness if we train on p	↑ ↓ ∞ ♥ ¥ added sequences and test on unpadded
C]	history = model.fit(train_dataset validation_da test_loss, test_act = model.eval print("Test Loss: [}'.format(test print("Test Accurecy: [}'.format) hove model does not mask the paddim	rt, epochs=2, iata=test_dataset) uate(test_dataset) t_loss)) (test_acc)) g applied to the sequ	ences. This can lead to skewness if we train on p	↑ ↓ ∞ ↓ added sequences and test on un padded

After setting the model; we will train the model for few epochs so that we get to experience output.

(Refer Slide Time: 29:46)



And we store the progress of the model in the history object, so that we can later plot how the training progressed.

(Refer Slide Time: 30:09)

ide + Text & Copy to Driv	•	RAM Disk
[] def pad_to_size(vec, si zeros = [0] * (size - vec.extend(zeros) return vec	ze): len(vec))	
<pre>[] def smple_predict(sent tokenized_sample_pred if pad: tokenized_sample_pr predictions = model.p return (predictions)</pre>	ence, pod); _text = tokenizer.encode(sample_pred_text) ed_text = pad_to_size(tokenized_sample_pred_text, 64) redict(tf.expand_diss(tokenized_sample_pred_text, 0))	
<pre>[] # predict on a simple t sample_pred_text = ('Th 'se predictions = sample_pr print (predictions)</pre>	ext without padding. a movie was cool. The animation and the graphics ' re out of this world. I would recommend this movie.') wdict(sample_pred_text, padwfaise)	
<pre>[] @ predict on a sample t sample_pred_text = ('Th 'we predictions = sample_pr print (predictions)</pre>	<pre>ext with padding e movie was cool. The animation and the graphics ' re out of this world. I would recommend this movie.') edict(sample_pred_text, pad=Frue)</pre>	
[] plot_graphs(history, 'a	coursey')	

The training usually takes longer to complete because we are trying to train on a large data set.

(Refer Slide Time: 30:31)



We can use the evaluate function on the model and obtain the test loss and the training loss. If the prediction has probability greater than 0.5, we mark it as a positive review; otherwise we mark it simply as a negative review.

We will also check how the model performance effects when we give sample text without padding and with padding. Ideally, the model should learn to ignore the padding, but you will experience that there is some effect of padding on the output.

(Refer Slide Time: 31:39)



So, now that you have trained our first model with a bi-directional LSTM, we will try to stack up couple of bi-directional LSTM and get model with more complexity. Let us see how to stack up different LSTMs and obtain in the model with more capacity. So, here we define the first bi-directional LSTM model; here we put return_sequence = true so that we get output from each node.

Each LSTM outputs 64 numbers in each direction, so the concatenation that happens in bi-directional LSTM will result into 128 numbers coming out of this particular layer. The second bi-directional LSTM will contribute to 64 numbers that will be passed into another dense layer followed by an output layer.

(Refer Slide Time: 32:55)



You can compile the model and fit the model and use the train model to calculate the loss and accuracy on the test set.

Just in case of the earlier model where we used a single bi-directional LSTM. This is the first example where we used bi-directional LSTM for predicting sentiments of movie reviews. We will have couple of more examples of using LSTM models for time series forecasting and for text generation.