

5-1 Sequence Model

5.1-1. Why choose sequence model

sequence models, one of the most exciting areas in deep learning. Models like recurrent neural networks or RNNs have transformed speech recognition, natural language processing and other areas

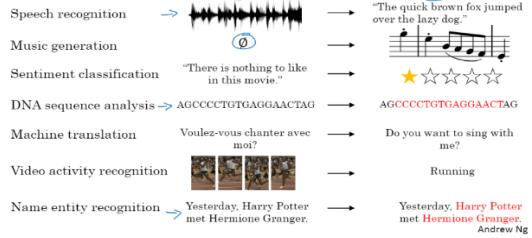
supervised learning with label data X as training set:

- CASES:
 1. input and output both are sequence, length maybe different;
 2. only one is sequence, the other not: Input/output

e.g. for using sequence:

1. Speech recognition: input audio.is sequence, output text is also sequence;
 Sequence models such as a recurrent neural networks and other variations, you'll learn about in a little bit have been very useful for speech recognition.
2. Music generation: only output special sequence. Input can be nothing or a signal trigger.
3. Sentiment classification: input 'human feeling feedback'-->star, sequence;
4. DNA sequence analysis: input DNA sequence, output sequence result
5. Machine translation: one sequence language to other language sequenced
6. Video activity recognition: sequence of pic -different gesture-->output activity result : like running
7. Name entity recognition: Give a sequence with name.-->output sequent name in the input

Examples of sequence data



5.1-2 Notions used in sequence model

1. Motivating example:

Name entity recognition: recognize name/country/location in different types of text, etc.
 used by search engine (搜索引擎) : like in the past 24h people name mentioned in news article

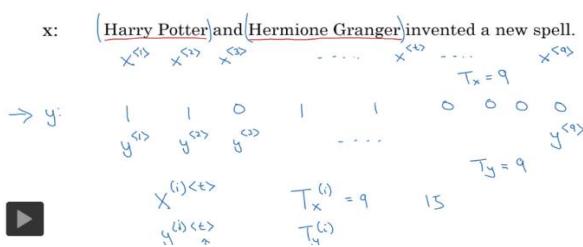
E.G: Named-entity recognition

input vector: sentence words

output vector: for each word whether it is part of a name: yes, 0 No.

Note: have more complicated output y, not only tell if there is a name, but also tell the start and end of the name.

Motivating example



1.1 Notes:

Input: $X = [X_{<1>}, X_{<2>}, \dots, X_{<t>}, \dots, X_{<Tx>}]$; T_x : length of input (number of input value);

Output: $Y = [y_{<1>}, y_{<2>}, \dots, y_{<t>}, \dots, y_{<Ty>}]$; T_y : length of output (number of output value)

t : index into the positions in the sequence.

input sample i: $X_{(i)}$.

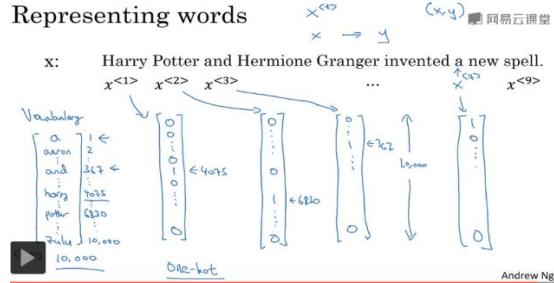
Training sample i, t position: $X_{(i)<t>}$

Output of sample i, t position: $y_{(i)<t>}$

Training samples can have different length, -->length of training sample i: $T_x(i)$

Output length of sample i: $T_y(i)$

Representing words



2. Representative for word:

>1. Come up a vocabulary /dictionary

A list of words that will use in word representations, commercial application may have 30-50,000 size/words dictionary.

Build dictionary: once chosen dictionary size like 10,000, and then could look through training sets, find the top 10,000 occurring words, or look through some of the online dictionaries that tells what are the most common 10,000 words in the English Language.

>2. Use one-hot vector to represent each of the input word: $x_{<t>}$

1: in the dictionary position corresponding to input $x_{<t>}$; other positions 0;

if no this word/value in dictionary, take a new fake word 'UNK' 'Unknown word' to represent it.

2. Representative for word:

Note:

1. one-hot representation of each word $x_{(i)<t>}$, is a vector; dimension is dictionary size.

2. Label $y_{(i)<t>}$: is also one-hot vector (output: one 1 and 0 is everywhere else)

3. The goal is given this representation for X to learn a mapping using a sequence model to then target output y . Take it as a supervised learning problem, with give labeled one-hot representation of both x and y .

NLP is supervised learning, with training data label (x, y), to train a neural network mapping x to y .

5.1-3 Recurrent Neural Networks

1. why not a standard network?

1.1. Standard network: $X[X_{<1>}, \dots, X_{<Tx>}]T \rightarrow \text{network} \rightarrow \text{output } Y[y_{<1>}, \dots, y_{<Ty>}]$

Note: each input word $x_{<t>}$ is one-hot vector

e.g:

input 9 words - each word is one-hot vector, feeding them into a standard neural network, maybe a few hidden layers, and then eventually had this output the nine values zero or one that tell whether each word is part of a person's name. But this turns out not to work well .

(note: a classifier for one word, output 1/0 - is a name or not, then go through each word of input.

not work as could not consider words dependency in input)

Standard network Problem:

1. Inputs and outputs length could be different in different input example ($x(i)$)

it's not as if every single example had the same input length T_x or the same output length, and even using padding/zero padding to make every input to that max length, this still does not seem a good representation.

2. Does not share features learned across different positions of text:

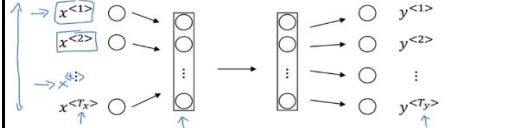
Features shared in different position means: if a value learned in input position $x_{<t>}$ is predicted as part of name, then automatically predict it as a part of name when the same value occurs in other positions in input $x_{<t>}$.-->kind of like

convert things learned in one part of image also generalize to other parts of the image (filter parameter same for all image positions), and we like a similar effects for sequence data as well.

3. high computation: W (features, Input : sum of all input $x(i)$: each position dimension * positions in input $x(i)$), enormous number of parameters.

e.g: input example's each feature $x(i)_t$, is dictionary size dimension vector (like dictionary size 1000, then $x(i)_t$ dimension is (1000,1)), input dimension : vocab size x word number

Why not a standard network?

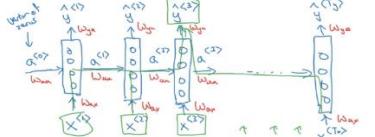


Problems:

Inputs, outputs can be different lengths in different examples.

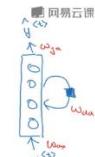
Doesn't share features learned across different positions of text.

Recurrent Neural Networks



He said, "Teddy Roosevelt was a great President."

He said, "Teddy bears are on sale!"



Andrew

2. Recurrent Neural Network (RNN)

RNN method without above standard networ problem: length difference, feature share positon-<---weight same in each position, --->less paramete, computation cost less

2.1 RNN idea:

eg: feed example i to neural network:
>>>1. Feed first word of example $x^{<1>} = \text{get output } y^{<1>} = a^{<1>}$
>>>2. Feed send word of example $x^{<2>} + a^{<1>} = \text{get output } y^{<2>} = a^{<2>}$, some info of first word computed, activation value from time step 1.

>>>3. Feed j_{th} word of example + previous word activation $a^{<j-1>} = \text{get output } y^{<j>} = a^{<j>}$

So at each time step, the recurrent neural network passes on previous time /word activation to the next time step for it to use

2.3 Advantage

1. Parameters are shared

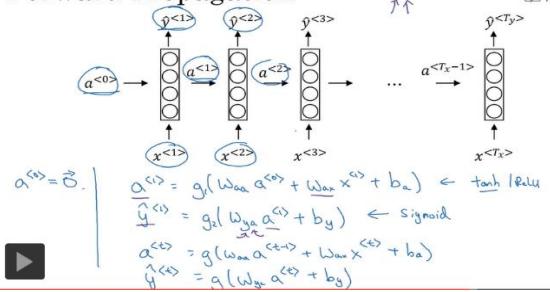
>>>1. Parameters used in each step are shared. Parameters: Weight on input $x^{<t>} : W_{ax}$, same for every time step/ all position $x^{<t>}$

>>>2. Horizontal activation function are governed by parameter W_{aa} : which is also same for every time step/ all position $x^{<t>}$

>>>3. Output predictions, govenred by W_{ya} : same for every time step/ all position $x^{<t>}$

>>>2. each time step output, will consider both current step input $x^{<t>}$, and info. in earlier time step: $x^{<1>} = \dots = x^{<t-1>}$. e.g. in time step 3, while making the prediction for $y^{<3>}$, algorithm gets the information not only from $x^{<3>}$ but also the information from $x^{<1>} \dots x^{<2>}$: because the information on $x^{<1>} \dots x^{<2>}$, can pass through the horizontal way/activation to help to prediction with $y^{<3>}$.

Forward Propagation



2.5 RNN Forward Propagation

>0. started off with the input $a^{<0>} = 0$ vectors.

Forward propagation:

>>>1. $a^{<1>} = g_a(W_{aa}a^{<0>} + W_{ax} * x^{<1>} + ba)$

>>>2. $y^{<1>} = g_y(W_{ya} * a^{<1>} + by)$

...

time step t:

$a^{<t>} = g_a(W_{aa}a^{<t-1>} + W_{ax} * x^{<t>} + ba)$

$y^{<t>} = g_y(W_{ya} * a^{<t>} + by)$

Activation function choice:

1. Activation function g_a : often is \tanh , sometimes are ReLU .

Have other ways of preventing the vanishing gradient problem.

(note: when to use ReLU , \tanh)

2. Activation function g_y : depending on what output y is.

If it is a binary classification problem, then a sigmoid activation function, or it could be a softmax that have a k-way classification problem.

The choice of activation function here will depend on what type of output y you have.

e.g.: name entity recognition task, output $y^{<t>} = 0$ or 1, g_y could be a sigmoid activation function.

5.1-4 Backpropagation through time

when implement this in one of the programming frameworks, often, the programming framework will automatically take care of backpropagation. But it's still useful to have a rough sense of how backprop works in RNNs.

1. Forwardprop and Backpropagation

1.1 Forwardprop

time step t:
 $a^{<t>} = g_a(W_{aa}a^{<t-1>} + W_{ax} * x^{(t)} + ba)$
 $y^{<t>} = g_y(W_{ya} * a^{<t>} + by)$

1.2: Backprop:

1. Loss function:

>>>1. Loss function with a signal prediciton at a singel position/time step $t^{<1>}$

e.g.: define an element-wise loss , which is supposed for a certain word in the sequence.
If it is a person's name, so labeled $y^{<t>} = 1$, neural network outputs some probability $y^{<t>} = p$, of the particular word being a person's name. So define lost function as the standard logistic regression loss,

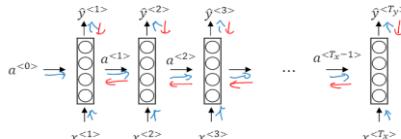
So this is the loss associated with a single prediction at a single position or at a single time set, t, for a single word.

>>>2. Total loss function L: sum over all the individual losses per timestep . (Ty number)

Loss for the entire sequence L will be defined as the sum over all t: 1, 2, ..., T_y. T_y is equals to T_y in this example. In a computation graph, to compute the loss: given $y^{<1>} \dots y^{<T_y>} \rightarrow$ compute the loss for the them independently: $<1> \dots <T_y>$, and then sum them all up to compute the final L.

Forward + Backprop is the computation graph.

Forward propagation and backpropagation



2. Recurrent Neural Network (RNN)

2.2 Note:

1. in this example length of input and output word are same $T_x = T_y$. Algorithm architecture need to modify a little bit if $T_x \neq T_y$.

2. to kick off the whole RNN, need made-up activation at time zero $a^{<0>} = 0$.
 $a^{<0>} = 0$ vector. Some researchers will initialized a 0 randomly. Could use other ways to initialize $a^{<0>} = 0$, but really having a 0 vector as the fake times 0's activation is the most common choice.

3. There are other type of neural network drawn-right pic, in which there is line to denote the recurrent connection, like a loop / shaded box to denote a time delay of one step. Left pic is the unrolled diagram, more easier to understand.

2.4: Disadvantage:

only use info. earlier in the sequence to make a prediction, do not consider info. in the later sequence.

While it is really helpful to know info. in both earlier and later sequence to make prediction for current positon. will solve in the BRNN(bidirectional recurrent neural network)

BRNN- Just have to make a quick modifications to these ideas later to enable, say, the prediction of $y^{<3>} = 0$ to use both information earlier in the sequence as well as information later in the sequence.

e.g: input "He said Teddy Roosevelt was a great president."

In order to decide whether or not the word Teddy is part of a person's name, it would be really useful to know not just information from the first two words but to know information from the later words in the sentence as well. because the sentence could also have been, "He said teddy bears are on sale." So given just the first three words is not possible to know for sure whether the word Teddy is part of a person's name.

Simplified RNN notation

Andrew Ng

2.6 Simplified RNN notation

for time step t:

$a^{<t>} = g_a(W_{aa}a^{<t-1>} + W_{ax} * x^{(t)} + ba)$

$y^{<t>} = g_y(W_{ya} * a^{<t>} + by)$

>1. combine W_{aa} , W_{ax} to matrix $W_a = [W_{aa} | W_{ax}]$, --> $a^{<t>} = g_a(W_a[a^{<t-1>} | x^{(t)}] + ba)$

Advantage: rather than carrying around two parameter matrices, W_{aa} and W_{ax} , we can compress them into just one parameter matrix W_a , and just to simplify our notation for when we develop more complex models.

Note:

1. Dimension: $a^{<t>}$ in different time step, have same dimension, = W_a . W_a row number.

e.g: $a[t]$ is $(n, 1)$, input $x^{<t>} = (n, d)$, d - dictionary size

--> $W_a: n \times n$, always square matrix; $W_{ax}: n \times n_d$

--> $W_a = [W_{aa} | W_{ax}]: n \times (n + n_d, 1)$

> 2. $[a^{<t-1>} | x^{(t)}] = [a^{<t-1>} | x^{(t)}]^T: (n + n_d, 1)$

>3. $W_y = W_{ya}$

2. Backpropagation through time:

Backprop then just requires doing computations in the opposite directions:

$da^{<t>} = dy^{<t>} * W_{ya} + g_y'(W_{ya} * a^{<t>} + by) + da^{<t+1>} * (-dy^{<t+1>} * W_{ya} + g_y'(W_{ya} * a^{<t+1>} + by)) * W_{ya} * g_a'$

>1. Vertical direction: $da^{<t>} = dy^{<t>} * W_{ya} * g_y'(W_{ya} * a^{<t>} + by)$

Total loss function --> each single loss function --> output $y^{<t>} \rightarrow$ activation $a^{<t>}$

>2. Horizontal direction: activation $a^{<t>} \rightarrow a^{<t-1>}$

Note:

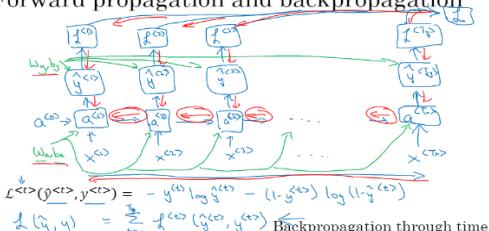
1. In this back propagation procedure, the most significant message or the most significant recursive calculation is Horizontal direction: activation $a^{<t>} \rightarrow a^{<t-1>}$

2. Horizontal back prop goes from right to left, that's reason this algorithm called backpropagation through time. The motivation for this name is that for forward prop, you are scanning from left to right, increasing indices of time t, whereas, the backpropagation, from right to left, kind of going backwards in time.

Summary:

So far, only seen this main motivating example in RNN in which the length of the input sequence was equal to the length of the output sequence.

Forward propagation and backpropagation



5.1-5 RNN archetecture type

Have seen an RNN architecture where the number of inputs, T_x , is equal to the number of outputs, T_y . It turns out that for other applications, T_x and T_y may not always be the same. It turns out that we could modify the basic RNN architecture to address all of these problems.

Here will see a much richer family of RNN architectures.

RNN architecture range

1: many to many: Input length = output length: $T_x = T_y$

2. Many to 1: only output AT the last time step when have the entire input sequence.

e.g. address sentiments classification

Input x might be a piece of text, such as it might be a movie review that says, "There is nothing to like in this movie." So x is going to be sequenced, and y might be a number from 1 to 5, or maybe 0 or 1.

In this case, we can simplify the neural network architecture as follows:

>1. Input the words one at a time.

If the input text was, "There is nothing to like in this movie." So "There is nothing to like in this movie," would be the input.

>2. Then rather than having to use an output at every single time-step, just have the RNN read into entire sentence and have it output y at the last time-step.

Examples of sequence data

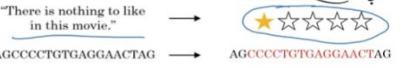
Speech recognition



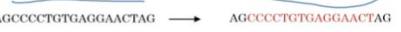
Music generation



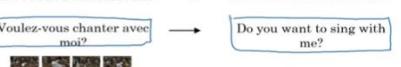
Sentiment classification



DNA sequence analysis



Machine translation



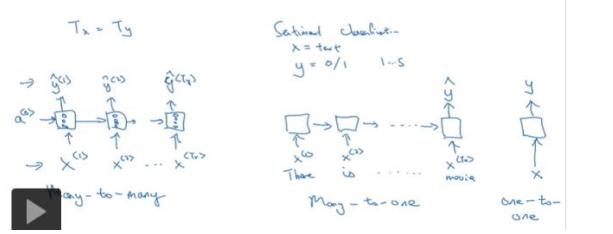
Video activity recognition



Name entity recognition



Examples of RNN architectures



RNN architecture range

3. One to many:

e.g. music generation:

output is a number of notes corresponding to a piece of music, input x could be an integer telling what kind of music want, or if don't want to input anything, x could be a null input, could always be the vector zeroes as well.

Neural network architecture:

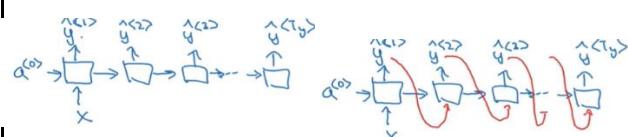
>1. Input $x, a^{<0>}$ have your RNN first output $y^{<t-1>}$

>2. With no further inputs, only activation $a^{<1>}$ in first step, output $y^{<2>}$

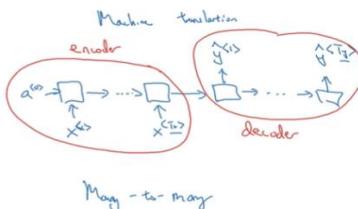
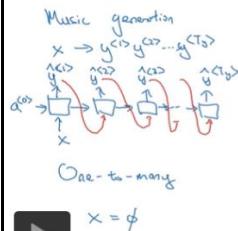
>...
> T_x : with no input X , only activation $a^{<Ty-1>}$, output $y^{<Ty>}$

Note: One technical used when you're actually generating sequences, often take the previous time step output $y^{<t-1>}$ as activation $a^{<t-1>}$ and feed it to the next layer.

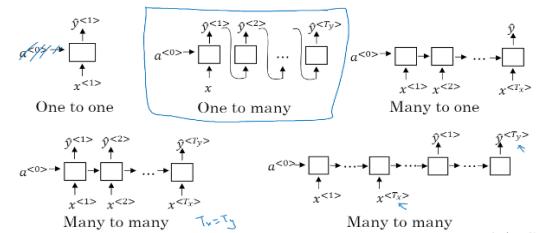
(then output $y^{<t>}$ is like one-one architecture, input is only $y^{<t-1>}$, output only depending on previous step output).



Examples of RNN architectures



Summary of RNN types



Andrew Ng

Summarize: the wide range of RNN architectures

there is one-to-many: like a music generation or sequenced generation;

there's many-to-one: like sentiment classification. Where you might want to read as input all the text with a movie review. And then, try to figure out that they liked the movie or not.

There is many-to-many: like the name entity recognition, where T_x is equal to T_y .

there's other version of many-to-many: like machine translation, T_x and T_y no longer have to be the same.

Now learned most of the building blocks, the building are pretty much all of these neural networks except that there are some subtleties with sequence generation, which is what we'll discuss in the next video.

5.1-6 Language Model and sequence generation

Language modeling is one of the most basic and important tasks in natural language processing. There's also one that RNNs do very well.

1. What is language modelling

1.1 Example: Input speech --> choose the best output text

Output choice 1: the apple and pair salad.

Output choice 2: the apple and pear salad.

You probably think the second sentence is much more likely, and in fact, that's what a good speech recognition system would help with even though these two sentences sound exactly the same.

The way a speech recognition system picks the second sentence is by using a language model which tells it what the probability is of either of these two sentences:

e.g. $P(\text{choice 1}) = 3.2 \times 10^{-13}$, $P(\text{choice 2}) = 5.7 \times 10^{-10}$.

With these probabilities, the second sentence is much more likely by over a factor of 10^{3} compared to the first sentence. And that's why speech recognition system will pick the second choice.

1. What is language modelling

1.2: Language model

>1. **Language model Function:** given any sentence, language model is to tell what is the probability of a sentence, of that particular sentence.

>2. **Probability of sentence:** means, if you want to pick up a random newspaper, open a random email or pick a random webpage or listen to the next thing someone says. What is the chance that the next sentence you use somewhere out there in the world will be a particular sentence like the apple and pear salad?

Note:

Language model is a fundamental component for both speech recognition systems, as well as for machine translation systems where translation systems wants output only sentences that are likely.

>3. **Language model basic job:**

Input a sentence, $y^{<1>} y^{<2>} \dots y^{<Ty>}$. language model estimates the probability of that particular sequence of words:

$P(y^{<1>} y^{<2>} \dots y^{<Ty>})$

What is language modelling?

Speech recognition

The apple and pair salad.

→ The apple and pear salad.

$$P(\text{The apple and pair salad}) = 3.2 \times 10^{-13}$$

$$P(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$$

$$P(\text{sentence}) = ?$$

$$P(y^{<1>} | y^{<2>} | \dots | y^{<t>})$$

2. Building language model with RNN

2.1 Build Id language model:

>1. first need Training set: comprising a large corpus of english/franch , or text from whatever language you want to build a language model.

Word corpus:NLP terminology means a large body /a very large set of english text, of english sentence.

E.g. one training sentence: 'Cat average 15 hours of sleep a day'

>2. Tokenize training sentence:

>>>1. build vocabulary/dictionary:

>>>2. Map the each of the word in training sentence to one-hot vector or the indices to vocabulary

>3. Commonly also add another extra token 'EOS', to the end of every sentence in the training set. EOS stands for sentence end and appended to the end of every sentence in the training set, so could help you figure out when a sentence ends. In this way training sentence length add one more: $y^{<T>} = \text{EOS}$.

3. RNN_Mode: I chance of difference sequences

e.g: one of training sentence is : cat average 15 hours of sleep a day. <EOS>.

RNN architecture:

>1. At time 0: make a soft max prediction to try to figure out what is the probability of the first words $y^{<1>} = P(\text{first word}?)$

computing activation $a^{<1>} = \text{softmax}(\text{some inputs} \times^{<1>} - \text{just set it to 0 vector: } x^{<1>} = 0 \text{ vector} + \text{and the previous } a^{<0>}), \text{ by convention, also set that to vector zeroes: } a^{<0>} = 0 \text{ vector.}$

$a^{<1>} = \text{softmax}(\text{some inputs} \times^{<1>} - \text{just set it to 0 vector: } x^{<1>} = 0 \text{ vector} + \text{and the previous } a^{<0>}), \text{ by convention, also set that to vector zeroes: } a^{<0>} = 0 \text{ vector.}$

Step time 0: get $y^{<1>} = \text{softmax}(\text{some inputs} \times^{<1>} - \text{just set it to 0 vector: } x^{<1>} = 0 \text{ vector} + \text{and the previous } a^{<0>}), \text{ by convention, also set that to vector zeroes: } a^{<0>} = 0 \text{ vector.}$

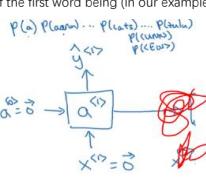
e.g: what's the chance that the first word is a, 'cats', 'Zulu', unknown word.

$y^{<1>} = \text{softmax}(\text{some inputs} \times^{<1>} - \text{just set it to 0 vector: } x^{<1>} = 0 \text{ vector} + \text{and the previous } a^{<0>}), \text{ by convention, also set that to vector zeroes: } a^{<0>} = 0 \text{ vector.}$

Note: Softmax function is trying to predict what is the probability of any word in the dictionary.

2. $y^{<1>} = \text{softmax}(\text{some inputs} \times^{<1>} - \text{just set it to 0 vector: } x^{<1>} = 0 \text{ vector} + \text{and the previous } a^{<0>}), \text{ by convention, also set that to vector zeroes: } a^{<0>} = 0 \text{ vector.}$

e.g: vocabulary have 10,000 words, then softmax have 10,000 output.



4. Train neural network: cost function definition

$$>1. \text{ Cost function at a certain time step } t, l(y^{<1>}, y^{<2>}, \dots, y^{<t>}) = -\sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

True word was $y^{<t>} = \text{one-hot vector, vocabulary size}$

, prediction $\hat{y}^{<t>} = \text{vector of vocabulary size.}$

--> loss at this time step (soft max loss function) : $y_i^{<t>} \log (\hat{y}_i^{<t>})$, sum over all vector size.

>2. Overall loss : sum over all time steps of the loss associated with the individual predictions.

Summarize:

And if you train this RNN on the last training set, what you'll be able to do is given any initial set of words, such as cats average 15 hours of sleep a day, can predict what is the chance of the next word.

e.g: new sentence: $y = y^{<1>} = y^{<2>} = y^{<3>} = \text{cats, average, 15, hours, of, sleep, a, day.}$ with just a three words, to get this three words sentence probability.

The way to figure out what is the chance of this entire sentence would be:

The first softmax tells he chance of $y^{<1>} = \text{cats}$. That would be this first output.

And then the second one can tell given the correct first work $y^{<1>} = \text{cats}$, what's the chance of $p(y^{<2>} | y^{<1>} = \text{cats})$

And then the third one tells what's the chance of $y^{<3>} = \text{average}$ given correct $y^{<1>} = \text{cats}$ & $y^{<2>} = \text{average}$: $p(y^{<3>} | y^{<1>} = \text{cats}, y^{<2>} = \text{average})$.

----> multiplying out these three probabilities, end up with the probability of the three-word sentence $p(y^{<1>} = \text{cats}, y^{<2>} = \text{average}, y^{<3>} = \text{sleep})$.

Sentence probability:

$$P(y^{<1>} = \text{cats}, y^{<2>} = \text{average}, y^{<3>} = \text{sleep}) = P(y^{<1>} = \text{cats}) * P(y^{<2>} = \text{average} | y^{<1>} = \text{cats}) * P(y^{<3>} = \text{sleep} | y^{<1>} = \text{cats}, y^{<2>} = \text{average})$$

So that's the basic structure of how you can train a language model using an RNN.
(note: language model is many-many RNN, $TxTy = Ty(Tx + 1.0 = Ty + (\text{EOS}))$; only output $y^{<t>} = \text{softmax}(a^{<t>})$ in each step is the chance of each vocab word being the output)

5.1-7 Sampling novel sequence_Recurrent neural network

After train a sequence model, one of the ways can informally get a sense of what is learned is to have a sample novel sequences.

1. Sampling a sequence from a trained RNN

A sequence model, models the chance of any particular sequence of words as follows, and so what we like to do is sample from this distribution to generate noble sequences of words.

To sample, need to do something slightly different:

>1. Time 0: sample the first word want model to generate.

>>>1. Inpute $x^{<1>} = 0, a^{<0>} = 0, \dots$ >>>output softmax probability , vector of vocabulary size.

>>>2. Then randomly sample according to this soft max distribution.

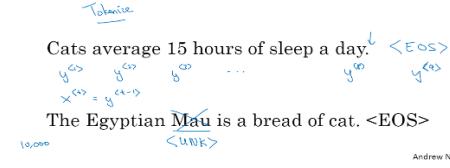
e.g: what the soft max distribution gives you is it tells you what is the chance that it refers to this a, what is the chance that it refers to this Aaron? What's the chance it refers to Zulu, what is the chance that the first word is the Unknown word token. Maybe it was a chance it was a end of sentence token.

And then take this vector $y^{<1>} = \text{softmax}(a^{<1>})$ and use like the numpy command `np.random.choice` to sample according to distribution defined by this vector probabilities, and that lets you sample the first words.

(note: could use a trained language model RNN to generate novel sequence)

Language modelling with an RNN

Training set: large corpus of english text.



Andrew Ng

2. Building language model with RNN

2.1 Build Id language model:

>1. first need Training set: comprising a large corpus of english/franch , or text from whatever language you want to build a language model.

Word corpus:NLP terminology means a large body /a very large set of english text, of english sentence.

E.g. one training sentence: 'Cat average 15 hours of sleep a day'

>2. Tokenize training sentence:

>>>1. build vocabulary/dictionary:

>>>2. Map the each of the word in training sentence to one-hot vector or the indices to vocabulary

>3. Commonly also add another extra token 'EOS', to the end of every sentence in the training set. EOS stands for sentence end and appended to the end of every sentence in the training set, so could help you figure out when a sentence ends. In this way training sentence length add one more: $y^{<T>} = \text{EOS}$.

3. RNN_Mode: I chance of difference sequences

e.g: one of training sentence is : cat average 15 hours of sleep a day. <EOS>.

RNN architecture:

>2. Time1: has activation $a^{<1>} = \text{try to figure out what is the second word: } p(\text{2nd word} = ? | \text{cats})$

Input: correct first word $y^{<1>} = x_2$,

RNN here use softmax functions to predict the chance of being whatever the second word it is, given what had come previously (correctly previous word).

Note: 1. softmax function also used;
2. $y^{<2>} = \text{vector of vocabulary size};$

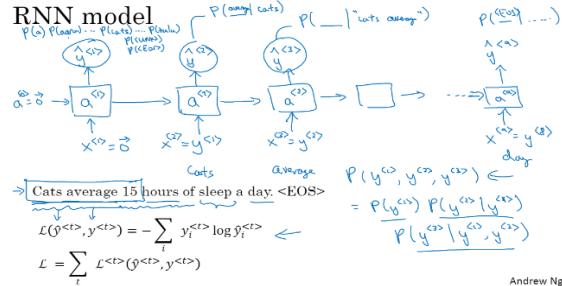
>3. Continue next step till the end of sentence, which is the EOS token.

Function: preface the chance of each vocabulary word, given all the correct word in earlier sequence: $p(\text{?} | \text{cat average 15 hours of sleep a day})$

Hopefully it will predict that's a high chance of it, EOS

So each step in the RNN will look at some set of preceding words such as, given the first three words, what is the distribution over the next word (chance of each vocabulary word being the i th word in sentence).

And so this RNN learns to predict one word at a time going from left to right.



Andrew Ng

5.1-7 Sampling novel sequence_Recurrent neural network

After train a sequence model, one of the ways can informally get a sense of what is learned is to have a sample novel sequences.

1. Sampling a sequence from a trained RNN

A sequence model, models the chance of any particular sequence of words as follows, and so what we like to do is sample from this distribution to generate noble sequences of words.

To sample, need to do something slightly different:

>1. Time 0: sample the first word want model to generate.

>>>1. Inpute $x^{<1>} = 0, a^{<0>} = 0, \dots$ >>>output softmax probability , vector of vocabulary size.

>>>2. Then randomly sample according to this soft max distribution.

e.g: what the soft max distribution gives you is it tells you what is the chance that it refers to this a, what is the chance that it refers to this Aaron? What's the chance it refers to Zulu, what is the chance that the first word is the Unknown word token. Maybe it was a chance it was a end of sentence token.

And then take this vector $y^{<1>} = \text{softmax}(a^{<1>})$ and use like the numpy command `np.random.choice` to sample according to distribution defined by this vector probabilities, and that lets you sample the first words.

(note: could use a trained language model RNN to generate novel sequence)

>2. Time1: sample $y^{<2>} = \text{softmax}(a^{<2>})$, given sampled $y^{<1>} = \text{softmax}(a^{<1>})$.

Input: instead using $y^{<1>} = \text{input}$, take the $y^{<1>} = \text{softmax}(a^{<1>})$ that just sampled, then this soft max will make a prediction for what is $y^{<2>} = \text{softmax}(a^{<2>})$.

e.g: sample the first word to be 'the', then you pass to time 1 as $x^{<2>} = y^{<1>} = \text{softmax}(a^{<1>})$.

And now you're trying to figure out what is the chance of what the second word is, given that the first word is 'the'. Then again use this type of sampling function to sample $y^{<2>} = \text{softmax}(a^{<2>})$.

>3. Go to next step: take whatever choice had sampled and pass that to next time step.

keep going until you get to the last time step.

Note:

1. how to know when the sequence ends.

>1. If the end of sentence token is part of vocabulary, then could keep sampling until generate an EOS token. that tells you have hit the end of a sentence and can stop.

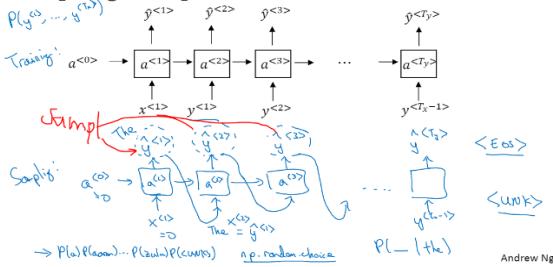
>2. Or alternatively, if do not include end token in vocabulary , then you can also just decide to sample 20 words or 100 words or something, and then keep going until you've reached that number of time steps.

2. This particular procedure will sometimes generate an unknown word token.

>1. If want to make sure that algorithm never generates this token, could just reject any sample that came out as unknown word token and just keep resampling from the rest of the vocabulary until you get a word that's not an unknown word.

>2. Or can just leave it in the output as well if you don't mind having an unknown word output.

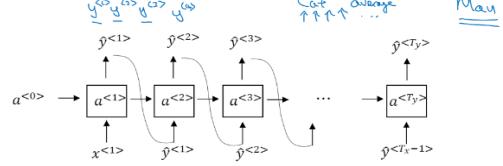
Sampling a sequence from a trained RNN



Character-level language model

→ Vocabulary = [a, aaron, ..., zulu, <UNK>]

→ Vocabulary = [a, b, c, ..., z, w, *, :, ;, 0, ..., a, A, ..., Z]



2.4: a character level language model pros and cons.

>Pros: Don't ever have to worry about unknown word tokens

In particular, a character level language model is able to assign a sequence like MAU, a non-zero probability. Whereas if MAU was not in vocabulary for the word level language model, you just have to assign it the unknown word token.

> Cons:

>>1. main disadvantage of the character level language model is that end up with much longer sequences.

So many English sentences will have 10 to 20 words but may have many, many dozens of characters. So character language models are not as good as word level language models at capturing long range dependencies between how the earlier parts of the sentence also affect the later part of the sentence.

>>2. character level models are also just more computationally expensive to train.

Character level model tends to be much harder, much more computationally expensive to train, so they are not widespread use today. Except for maybe specialized applications where you might need to deal with unknown words or other vocabulary words a lot. Or they are also used in more specialized applications where you have a more specialized vocabulary.

The trend in natural language processing is that for the most part, word level language model are still used, but as computers gets faster there are more and more applications where people are, at least in some special cases, starting to look at more character level models.

So under these methods, can build an RNN to look at the purpose of English text, build a word level, build a character language model, sample from the language model that have trained.

Sequence generation

3. Sequence generation example

E.g1_News

Text generated from a character-level language model. The text looks vaguely like news text, not quite grammatical, but maybe sounds a little bit like things that could be appearing news.

E.g2_Shakespeare

Language model was trained on Shakespearean text, and then it generates stuff that sounds like Shakespeare could have written it. The mortal.

Summary:

So this is it for the basic RNN, and how you can build a language model using it, as well as sample from the language model that you've trained.

(basic RNN --> Language model --> sampling from Language model: generating noval sentence

Language model: modified on basic RNN, input x<t>= correct y<t-1>

>>>softmax predict the next word given previous words: p(y|t+1 | y<1>, ..., y<t>) (note: loss is: y<t>-one vector * prediction y^<t>)

>>>for generating y^<t>: input correct y<t-1>=x<t>, and a<t-1>; while y^<1> except, input x<1>=0, a<0>=0

Sampling/generating noval sentece: use a trained language model, while setting x<t>= sampled y^<t-1>

In the next few videos, I want to discuss further some of the challenges of training RNNs, as well as how to adjust some of these challenges, specifically vanishing gradients by building even more powerful models of the RNN. So in the next video let's talk about the problem of vanishing the gradient and we will go on to talk about the GRU, Gate Recurring Unit as well as the LSTM models.

5.1-8 Vanishing gradients with RNNs

You've learned about how RNNs work and how they can be applied to problems like name entity recognition, as well as to language modeling, and you saw how backpropagation can be used to train in RNN. It turns out that one of the problems with a basic RNN algorithm is that it runs into vanishing gradient problems.

1. Vanishing gradients with RNNs

Basic RNN run into vanishing gradient problems.

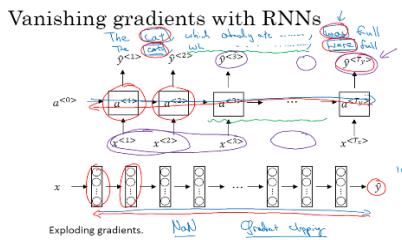
1.1 Language long-term dependency:

Language can have very long-term dependencies, where it worked at this much earlier can affect what needs to come much later in the sentence.

E.g: below two sentence, cat or cats can affect to use was or were.

The cat, which already ate....., was full;

The 'cats', which already ate...., 'were' full.



1. Vanishing gradients with RNNs

Basic RNN run into vanishing gradient problems.

1.2: Basics RNN is not very good at capturing very long-term dependencies.

>1. Vanishing gradient problem in deep neural network:

In a very, very deep neural network like 100 layers or even much deeper, then you would carry out forward prop, from left to right and then back prop. And in a deep neural network, the gradient from output y, would have a very hard time propagating back to affect the weights of these earlier layers, to affect the computations in the earlier layers.

>2. RNN Vanishing gradient problem:

RNN has a similar problem: have forward prop came from left to right, and then back prop, going from right to left. And it can be quite difficult, because of the same vanishing gradients problem, for the outputs, of the errors associated with the later time steps to affect the computations that are earlier.

>>>And so in practice, it means: it might be difficult to get a neural network to realize that it needs to memorize the just see a singular noun or a plural noun, so that later on in the sequence that can generate either was or were, depending on whether it was singular or plural.

And in English, this stuff in the middle could be arbitrarily long, then might need to memorize the singular/plural for a very long time before you get to use that bit of information.

1.2: Basics RNN is not very good at capturing very long-term dependencies

Summarize:

---> due to vanishing gradient, basic RNN model has many local influences: , meaning output like y^<3> is mainly influenced by values close to y^<3>, and it's difficult to be strongly influenced by an input that was very early in the sequence.

And this is because whatever the output is, whether this got it right or wrong, it's just very difficult for the area to backpropagate all the way to the beginning of the sequence, and therefore to modify how the neural network is doing computations earlier in the sequence.

So this is a weakness of the basic RNN algorithm.

2 Language model exploring gradient problem

Neural network exploring gradients:

>1. Also have exploring gradient problem

when doing back prop, the gradients should not just decrease exponentially, they may also increase exponentially with the number of layers you go through.

>2. Exploring gradient problem easy to find:

It turns out that vanishing gradients tends to be the bigger problem with training RNNs, although when exploding gradients happens, it can be catastrophic because the exponentially large gradients can cause your parameters to become so large that your neural network parameters get really messed up, so easier to spot because the parameters just blow up and you might often see NaNs, or not a numbers, meaning results of a numerical overflow in your neural network computation.

>3. Exploring gradient solution: apply gradient clipping.

Gradient clipping: if gradient vector is bigger than some threshold, re-scale some of your gradient vector so that is not too big.

Clipping method: could clip according to some maximum value.

(note: reduce gradient step by clipping, direction no change)

Summarize:

So there is exploding gradients, the derivatives do explode or you see NaNs, then just apply gradient clipping, and that's a relatively robust solution that will take care of exploding gradients.

But vanishing gradients is much harder to solve and it will be the subject of the next few videos.

5.1-9 Gate Recurrent Unit (GRU)

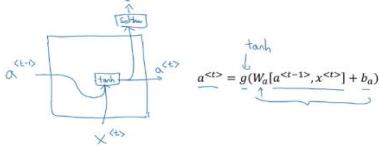
Gated Recurrent Unit is a modification to the RNN hidden layer that makes it much better capturing long range connections and helps a lot with the vanishing gradient problems.

1. RNN unit

Visualization of the RNN unit of the hidden layer: right pic.

$a^{t-1} \rightarrow \text{softmax} \rightarrow \text{output } y^{t-1}$.

RNN unit



2. GRU (Simplified) (modified on RNN):

Purpose: capturing long term dependency, prevent gradient vanishing

2.1. New cell: memory cell: c^{t-1}

Like previous sentence example, need to remember the cat was singular, to make sure you understand why that was rather than were. So the cat-was r or the cats-were.

Memory cell function: c

So as we read in this sentence from left to right, the GRU unit is going to have a new variable called c. stands for memory cell, to provide a bit of memory to remember, like, whether cat was singular or plural, so that when it gets much further into the sentence it can still work under consideration whether the subject of the sentence was singular or plural.

Note:

1. in GRU $c^{t-1} = a^{t-1}$ (in LSTM c^{t-1} and a^{t-1} is different value)

2.2. Candidate for replacing c^{t-1} : $C^{t-1} = \tanh(W_c [c^{t-1}, x^{t-1}] + b_c)$

At every time-step, we're going to consider overwriting the memory cell c^{t-1} with a value c^{t-1}

2.3. Gate Gu: (between 0-1): sigmoid ($W_u [c^{t-1}, x^{t-1}] + b_u$):

Gu will be a value between(0,1), to develop intuition about how GRUs work, think it as being always 0 or 1.

In practice, compute Gu with a sigmoid function, and it's value is always between (0,1). And for most of the possible ranges of the input (too large or small), the sigmoid function is either very, very close to zero or very, very close to one. So for intuition, think of gamma Gu as being either 0 or 1 most of the time.

2.5 GRU unit picture

Input: x^{t-1} , $c^{t-1} = a^{t-1}$

Output:

$C^{t-1} = \tanh(W_c [c^{t-1}, x^{t-1}] + b_c)$

$Gu = \text{sigmoid} ((W_u [c^{t-1}, x^{t-1}] + b_u))$

$C^{t-1} = Gu * C^{t-1} + (1-Gu) * C^{t-1}$

$y^{t-1} = \text{softmax} (W_y * C^{t-1} + b_y)$

GRU unit takes x^{t-1} , $c^{t-1} = a^{t-1}$ as input and together x^{t-1} and $c^{t-1} = a^{t-1}$ are passed through a tanh function to produce the new value for the memory cell $c^{t-1} = a^{t-1}$. which could also use softmax or something to make some prediction for y^{t-1} .

Above is he GRU unit or at least a slightly simplified version of it.

And what is remarkably good at is through the gates deciding that when you're scanning the sentence from left to right say, that's a good time to update one particular memory cell and then to not change, not change it until you get to the point where you really need it to use this memory cell that is set even earlier in the sentence.

2.6. Help vanishing gradient:

Because gamma Gu can be so close to 0, then it doesn't suffer from much of a vanishing gradient problem:

The gate Gu is quite easy to set to 0, as long as sigmoid content is large negative. In this case, $c^{t-1} = a^{t-1}$: this is very good at maintaining the value for the cell. c^{t-1} is maintained pretty much exactly even across many many many time steps.

So this can help significantly with the vanishing gradient problem and therefore allow a neural network to go on even very long range dependencies, such as a cat and was related even if they're separated by a lot of words in the middle.

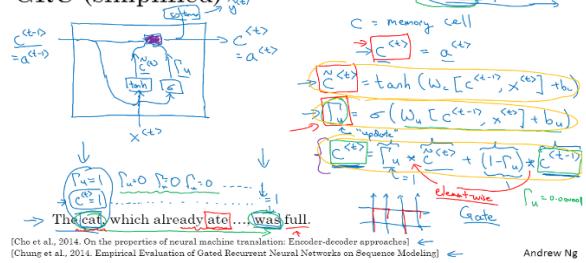
(note: backward prop: $dc^{t-1} = f(dc^{t+1})$; nearly: $dc^{t-1} \sim dc^{t+1}$ due to sigmoid and tanh function derivative: $s' = s(1-s)$.)

3. Summary:

While training of very deep neural network, you can run into a vanishing gradient or exploding gradient problems with the derivative, either decreases exponentially or grows exponentially as a function of the number of layers. And in RNN processing data over like a thousand times sets, over 10,000 times sets, that's basically a 1,000 layer or 10,000 layer neural network, and so, it too runs into these types of problems. Exploding gradients could sort of address by just using gradient clipping, but vanishing gradients will take more work to address.

So what we do in the next video is talk about GRU, the gate recurrent units, which is a very effective solution for addressing the vanishing gradient problem and will allow your neural network to capture much longer range dependencies.

GRU (simplified)



Andrew Ng

2.4 Key part of GRU: memory cell will be updated by candidate c^{t-1} , and gate Gu will decide whether or not actually update c with $c \leftarrow c^{t-1} = Gu * c^{t-1} + (1-Gu) * c^{t-1}$

We have come up with a candidate C^{t-1} , where we're thinking of updating c^{t-1} , and the gate Gu will decide whether or not we actually update it.

e.g. The cat, which already ate... was full:

And so the way to think about it is:

>1. In 'cat' time step, memory cell C^{t-1} is set to either 0 or 1 depending on whether the word you are considering, really the subject of the sentence is singular or plural (e.g. if sentence subject cat is singular, set $c^{t-1}=1$, otherwise set $c^{t-1}=0$).

>2. then the GRU unit would memorize the value of the c^{t-1} all the way until 'was' time step, where c^{t-1} still = 1, tells sentence object is singular so use the choice 'was'.

>1. The job of the gate Gu , is to decide when to update c^{t-1} :

e.g. when reach 'cat' time step, realize that is sentence object and update $c^{t-1}=1$, then go to next time step, no need to update c^{t-1} , and thus set Gate $Gu=0$. Until reach 'was' time step, use c^{t-1} to choose 'was'. then don't need to memorize anymore, I can just forget that.

>2. Gate formulation: $c^{t-1} = Gu * c^{t-1} + (1-Gu) * c^{t-1}$

If $Gu = 1$, new value of $c^{t-1} = c^{t-1}$

e.g. $Gu = 1$ in 'was' time step, then for all of these values in the middle, $Gu=0$: do not update c^{t-1} , just hang onto the old value, till 'was' time step: c^{t-1} it still memorizes, the cat was singular.

If $Gu=0$, new value of $c^{t-1} = c^{t-1}$ old value.

>2.7 : matrix dimension:

While implementing GRU:

>1. C^{t-1} can be a vector:

e.g. if have 100 dimensional or hidden activation value then c^{t-1} can be a 100 dimensional.

>2. C^{t-1} , Gu : are the same dimension with C^{t-1}

>3. Formal: $C^{t-1} = Gu * C^{t-1} + (1-Gu) * C^{t-1}$

e.g. if gate Gu is 100 dimensional vector, what it is really a 100 dimensional vector of bits, the value is mostly zero and one. That tells of this 100 dimensional memory cell which are the bits you want to update.

Note:

1. in practice gate Gu won't be exactly zero or one. Sometimes it takes values in the middle as well but it is convenient for intuition to think of it as mostly taking on values that are exactly zero, or pretty much exactly one.

2. these element wise multiplications just tells GRU which are the dimensions of memory cell vector to update at every time-step.

So you can choose to keep some bits constant while updating other bits:

e.g. maybe you use one bit to remember the singular or plural cat and maybe use some other bits to realize that you're talking about food.

You can use different bits and change only a subset of the bits every point in time.

Here is actually a slightly simplified GRU unit.

Full GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\tilde{c}_{t-1}^{<t-1>}, x^{<t>}] + b_c)$$

$$u \left\{ \begin{array}{l} \Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r) \end{array} \right.$$

$$h \quad c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

LSTM

The cat, which ate already, was full.

3. Full GRU unit

only one change, that is for $c^{<t>}$: $C^{<t>} = \tanh(W_c[\mathbf{Gr} \cdot c^{<t-1>}, x^{<t>}] + bc)$

One more gate \mathbf{Gr} added: $\mathbf{Gr} = \text{sigmoid}(W_r * [c^{<t-1>}, x^{<t>}] + bc)$
 \mathbf{Gr} tells how relevant is $c^{<t-1>}$ to computing the next candidate for $c^{<t>}$.

Summary:
 There are multiple ways to design these types of neural networks. And why do we have gate \mathbf{Gr} . Why not use a simpler GRU. So it turns out that over many years researchers have experimented with many, different possible versions of how to design these units, to try to have longer range connections, to try to have more of the longer range effects and also address vanishing gradient problems.

And the GRU is one of the most commonly used versions that researchers have converged to and found as robust and useful for many different problems.

And the other common version is called an LSTM which stands for Long Short Term Memory. GRUs and LSTMs are two specific instantiations of this set of ideas that are most commonly used.

Note: there could be different notation used at academic literature, like h instead of $c^{<t>}$, use $h^{<t>}$ etc.

So that's it for the GRU, for the Gate Recurrent Unit. This is one of the ideas in RNN that has enabled them to become much better at capturing very long range dependencies has made RNN much more effective. the other most commonly used variation of this class of idea is something called the LSTM unit, Long Short Term Memory unit.

5.1.10 LSTM (Long short term memory) Unit

The gated recurrent units, and how that can allow you to learn very long range connections in a sequence. The other type of unit that allows you to do this very well is the LSTM or the long short term memory units, and this is even more powerful than the GRU.

1. GRU and LSTM

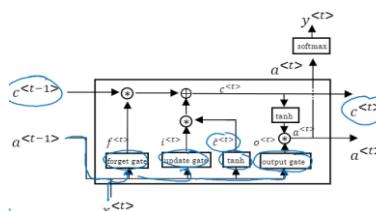
The LSTM is an even slightly more powerful and more general version of the GRU, and is due to Sepp Hochreiter and Jürgen Schmidhuber. And this was a really seminal paper, a huge impact on sequence modelling and Deep Learning community.

1.1 LSTM Compare vs GRU

- >1. $c^{<t>}$ not equal to $a^{<t>}$.
- >2. candidate for replacing $c^{<t>}$: $C^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + bc)$, no relevant gate \mathbf{Gr}
- >3. Update Gate \mathbf{Gu} is same;
- >4. Have one more gate: **forget gate \mathbf{Gf}** to control $c^{<t>}$: $Gf = \text{sigmoid}(W_f[a^{<t-1>}, x^{<t>}] + bf)$
- >5. Have another one more gate: **output gate \mathbf{Go}** = sigmoid ($W_o[a^{<t-1>}, x^{<t>}] + bo$);
- >6. $c^{<t>} = Gu * c^{<t-1>} + Gf * c^{<t-1>}$
 This is a vector-vector element-wise multiplication, have a separate forget gate Gf : gives the memory cell the option of keeping the old value $c^{<t-1>}$ and then just adding to $c^{<t>}$. So, use a separate update and forget gates.
- >7. $a^{<t>} = Go * \tanh(c^{<t>})$ (why make $a^{<t>} \neq c^{<t>}$)
- >8. $y^{<t>} = \text{softmax}(a^{<t>})$

these equations govern LSTM behavior

2. LSTM unit and network picture



simplified to diagrams a little bit in the bottom, like right pic:

- >1. $C<0> \dots >c<t> \dots >c<T>$
- >2. $a<0> \dots >a<t> \dots >a<T>$

the cool thing is: for the top horizontal line $c<0> \dots c<t>$, so long as set the forget Gf and the update gate Gu appropriately, it is relatively easy for the LSTM to have some value $c<0>$ and have that be passed all the way to the right to have maybe, $c<3> = c<0>$.

And this is why the LSTM, as well as the GRU, is very good at memorizing certain values even for a long time, for certain real values stored in the memory cell even for many, many timesteps.

(note: all the gate decided by activation and input: $a^{<t-1>}, x^{<t>}$, while $a^{<t-1>} = f(c^{<t-1>})$
 $\dots \rightarrow$ gate actually decided by: $c^{<t>}, x^{<t>}$)

3. Other version LSTM

The most common one is that instead of just having the gate values be dependent only on $a^{<t-1>}$, & $x^{<t>}$, sometimes, people also sneak in the values $c^{<t-1>}$ as well. This is called a peephole connection.

Peephole connection: means the gate values may depend not just on $a^{<t-1>}$, & $x^{<t>}$, but also on the previous memory cell value $c^{<t-1>}$.
 And the peephole connection can go into all three of these gates' computations.
 $Gu = \text{sigmoid}(W_u[a^{<t-1>}, x^{<t>}, c^{<t-1>} + bu])$
 $Gf = \text{sigmoid}(W_f[a^{<t-1>}, x^{<t>}, c^{<t-1>} + bf])$
 $Go = \text{sigmoid}(W_o[a^{<t-1>}, x^{<t>}, c^{<t-1>} + bo])$

Note:
 One technical detail is that these are, like 100-dimensional vectors. So if you have a 100-dimensional hidden memory cell unit, and the fifth element of $c^{<t-1>}$ affects only the fifth element of the corresponding gates, so that relationship ($c^{<t-1>}$, and gates) is one-to-one, where not every element of the 100-dimensional $c^{<t-1>}$ can affect all elements of the case. But instead, the first element of $c^{<t-1>}$ affects the first element of the case, second element affects the second element, and so on.

If you ever read the paper and see someone talk about the peephole connection, that's when they mean that $c^{<t-1>}$ is used to affect the gate value as well.

LSTM in pictures

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

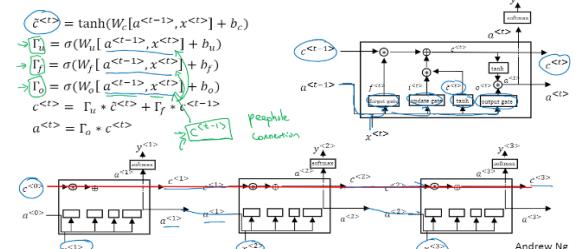
$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh(c^{<t>})$$

$$\tilde{c}^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh(c^{<t>})$$

LSTM in pictures



5.1.11: Bidirectional RNN (BRNN)

By now, you've seen most of the cheap building blocks of RNNs.
 But, there are just two more ideas that let you build much more powerful models:

One is bidirectional RNNs, which lets you at a point in time to take information from both earlier and later in the sequence, so we'll talk about that in this video.
 And second, is deep RNNs.

Summary:

>1. When to choose GRU or LSTM:
 There isn't widespread rule in this. LSTMs actually came much earlier, and then GRUs were relatively recent invention that were maybe derived as Pavia's simplification of the more complicated LSTM model.

Researchers have tried both of these models on many different problems, and on different problems, different algorithms will win out. So, there isn't a universally-superior algorithm which is why I want to show you both of them

Advantage of the GRU: it's a simpler model and so it is actually easier to build a much bigger network, it only has two gates, so computationally, it runs a bit faster. So, it scales the building somewhat bigger models.

LSTM: is more powerful and more effective since it has three gates instead of two.

If need to choose one, LSTM has been the historically more proven choice, most people today will still use the LSTM as the default first thing to try.

In the last few years, GRUs had been gaining a lot of momentum and I feel like more and more teams are also using GRUs because they're a bit simpler but often work just as well. It might be easier to scale them to even bigger problems.

1. Getting information from the future

1.1 Unidirectional or forward directional only RNN :

e.g:

He said "Teddy bears are on sale!"

He said " Teddy Roosevelt was a great President!"

-> To recognize name in sentence.

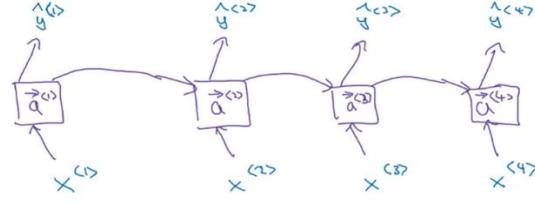
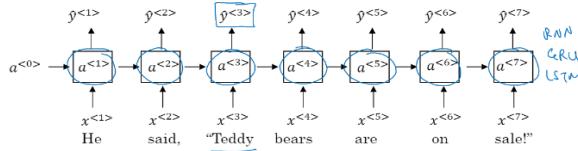
One of the problems of this network is that, to figure out whether the third word Teddy is a part of the person's name, it's not enough to just look at the first part of the sentence, which doesn't tell you if they're talking about Teddy bears or talk about the former US president, Teddy Roosevelt.

So this is a unidirectional or forward directional only RNN., whether these cells are standard RNN blocks or GRU units or LSTM blocks. all of these blocks are in a forward only direction.

A bidirectional RNN or BRNN, could fix this issue.

Getting information from the future

He said, "Teddy bears are on sale!"
He said, "Teddy Roosevelt was a great President!"



2. Bidirectional RNN (BRNN)

e.g:

Input: $x<1>, x<2>, x<3>, x<4>$
activation in forward: $a<1>, a<2>, a<3>, a<4>$
activation in backward: $a<4>, a<3>, a<2>, a<1>$
output: $y<1>, y<2>, y<3>, y<4>$

2.1: Forward prop: same as basic RNN:

Horizontally: $a<t-1>, x<t> \dots \rightarrow a<t>$

vertically: $a<t> (f(x<t>, a<t-1>)) \rightarrow y<t>$

2.2: Backward prop:

Horizontally: $a<t>, x<t-1> \dots \rightarrow a<t-1>, b$

vertically: $a<t-1>, b \rightarrow y<t>$

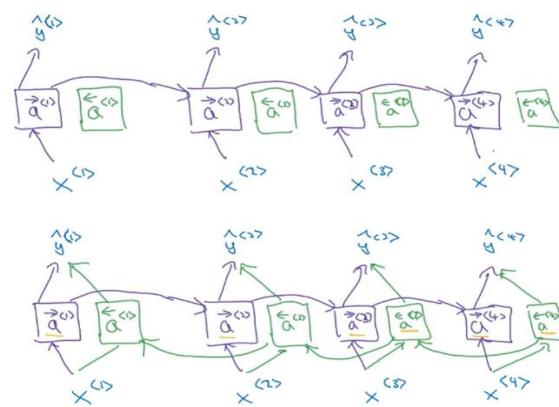
And this a \leftarrow activation backward connections will be connected to each other going backward in time.

the forward and backward network defines a a cyclic graph (无环图).

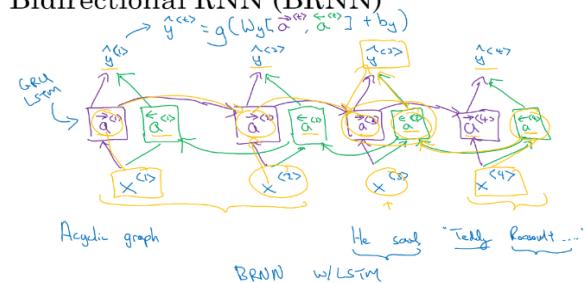
Given an input sequence $x<1> \dots x<4>$, forward process first compute $a<1>$ forward , then to $a<2>$ forward, then to $a<3>$ forward, then $a<4>$ forward . Whereas, the backward sequence would start by computing $a<4>, b$ backward , and then to $a<3>, b$ backward, ...to $a<1>, b$ backward.

Note: as computing network activation in backward, this is not backward actually is forward prop. But the forward prop has part of the computation going from left to right and backward prop has part of computation going from right to left in this diagram.

(note: parameter wa,ba, should be different in forward and backward?)



Bidirectional RNN (BRNN)



2.3: Make predictions $y<t>$:

after have computed all the activations in forward and backward, can make predictions, based on both forward $a<t>$ and backward activation $a<t>, b$:
 $y<t> = g(Wy[a<t>, a<t>, b] + by)$

e.g: for time step 3:

Input from current

>1. $x<3> \rightarrow a<3>$

>2. $x<3> \rightarrow a<3>, b$

current info. goes into both forward and backward.

Input from forward:

>1. $x<1> \rightarrow a<1>, \dots \rightarrow a<2>, \dots \rightarrow a<3>$

>2. $x<2> \rightarrow a<2>, \dots \rightarrow a<3>$

Input from backward:

>1. $x<4> \rightarrow a<4>, \dots \rightarrow a<3>, b$

Output: $y<3> = y<t> = g(Wy[a<3>, a<3>, b] + by)$

Forward and backward allows the prediction at time three to take as input both information from: 1. the past, 2. as well as information from the present which goes into both the forward and the backward things at this step.3. as well as information from the future.

Summary:

1. BRNN block:

The bidirectional recurrent neural network blocks here (represented by activation function $a<t>, a<t>, b$), can be standard RNN bloc or GRU blocks or LSTM blocks.
In fact, for a lots of NLP problems, a bidirectional RNN with a LSTM appears to be commonly used.

So, we have NLP problem and you have the complete sentence, you try to label things in the sentence, a bidirectional RNN with LSTM blocks both forward and backward would be a pretty views of first thing to try.

2. BRNN Advantage:

BRNN is able to make predictions anywhere even in the middle of a sequence by taking into account information potentially from the entire sequence.

3. BRNN disadvantage : do need the entire sequence of data before can make predictions anywhere.

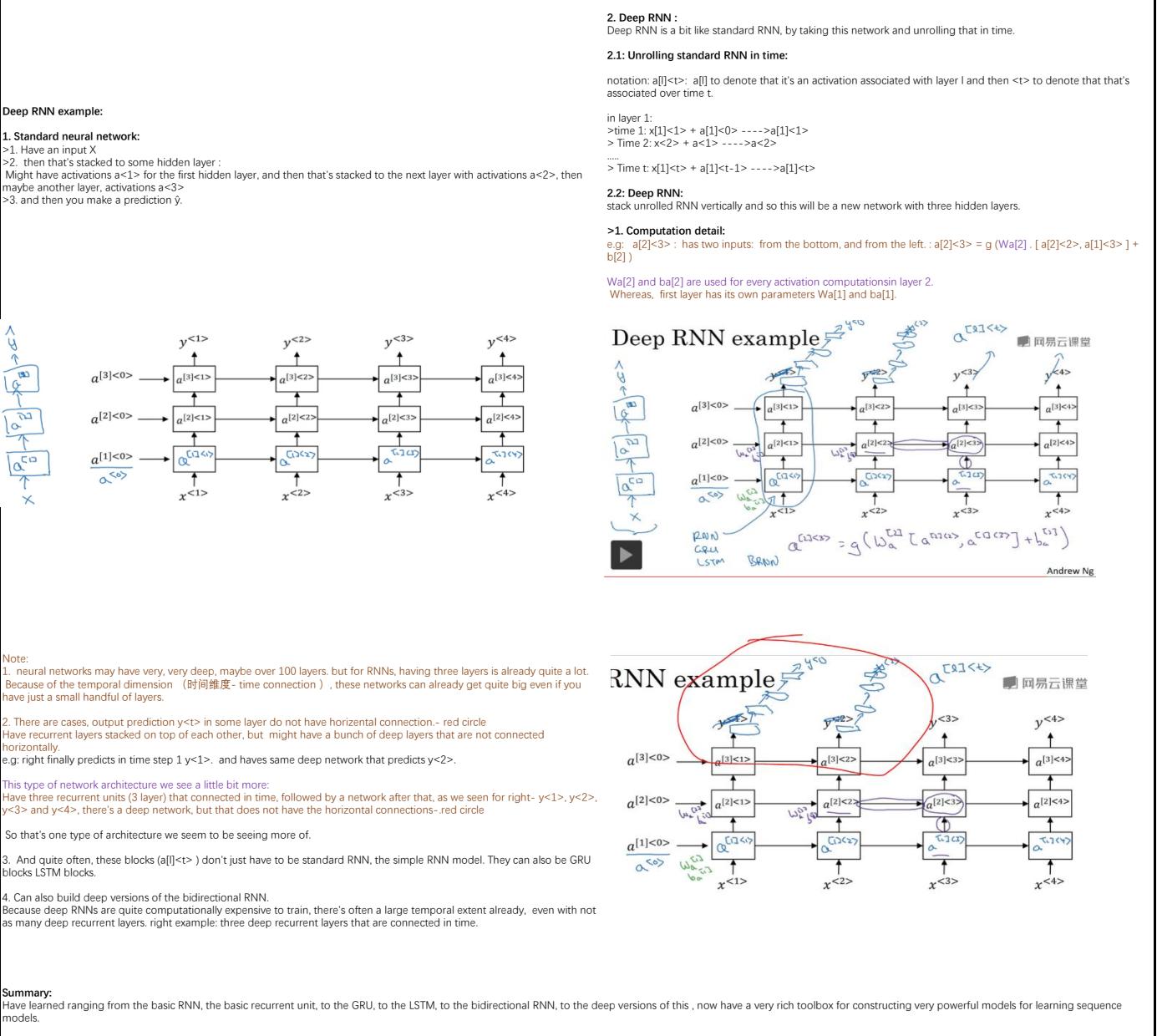
e.g: while building a speech recognition system, BRNN will let you take into account the entire speech utterance. using BRR directly, need to wait for the person to stop talking to get the entire utterance before you can actually process it and make a speech recognition prediction.

So for a real type speech recognition applications, they're somewhat more complex modules as well rather than just using the standard bidirectional RNN showed here.

But for a lot of natural language processing applications where you can get the entire sentence all the same time, the standard BRNN algorithm is actually very effective.

5.1-12 Deep RNNs

The different versions of RNNs you've seen so far will already work quite well by themselves. But for learning very complex functions sometimes is useful to stack multiple layers of RNNs together to build even deeper versions of these models.



5-2_NLP and Word Embedding

5.2-1 Word representation_NLP and word embedding

Last week, we learned about RNNs, GRUs, and LSTMs. In this week, you see how many of these ideas can be applied to NLP, to Natural Language Processing, which is one of the features of AI because it's really being revolutionized by deep learning.

One of the key ideas is word embeddings: which is a way of representing words. That let your algorithms automatically understand analogies like that, man is to woman, as king is to queen, and many other examples.

And through these ideas of word embeddings, you'll be able to build NLP applications, even with models the size of relatively small label training sets.

Finally towards the end of the week, you'll see how to debias word embeddings: that's to reduce undesirable gender or ethnicity or other types of bias that learning algorithms can sometimes pick up.

1. Word representation:

1.1 1-hot representation

We've been representing words using a vocabulary of words, and we've been representing words using a one-hot vector.

> 1-hot vector:

e.g., if man is word number 5391 in this dictionary, then represent it with a vector with 1 in position 5391, use \mathbf{O}_{5391} to represent this vector, O here stands for one-hot. all the word will be similarly represented with one-hot vector.

>2 one weaknesses of 1-hot representation:

It treats each word as a thing onto itself, and it doesn't allow an algorithm to easily generalize the cross words:

e.g.: have a language model that has learned that when you see sentence 'I want a glass of orange juice'. the blank word will be very likely 'juice'.

But even if the learning algorithm has learned that "I want a glass of orange juice" is a likely sentence", if it sees "I want a glass of apple". As far as it knows the relationship between apple and orange is not any closer as the relationship between any of the other words man, woman, king, queen, and orange,

It's not easy for the learning algorithm to generalize from knowing that orange juice is a popular thing, to recognizing that apple juice might also be a popular thing or a popular phrase. And this is because the any product between any two different one-hot vector is zero.

e.g: If take any two vectors like queen and king and any product of them, the end product is zero. If take apple and orange and any product of them, the end product is also zero.

--> Couldn't get the distance between any pair of these vectors, as all are same 0. So it just doesn't know that somehow apple and orange are much more similar than king and orange or queen and orange.

So, it will be nice if instead of a one-hot presentation we can instead learn a featurized representation with each of these words.

1.2 featurized representation

>1. Feature Intuition:

Featurized representation with each of these words: A man, woman, king, queen, apple, orange or really for every word in the dictionary, we could learn a set of features and values for each of them.

>2. for each of right words, we want to know:

>>1. Feature-Gender: what is the gender associated with each of these words.

If gender goes from -1 for male to +1 for female, then the gender associated with man might be -1, for woman might be +1. And then eventually, learning these things maybe for king you get -0.95, for queen+ 0.97, and for apple and orange sort of genderless.

>>2. Feature -Royal: how royal are these word.

Words man and woman are not really royal, so they might have feature values close to 0. Whereas king and queen are highly royal. And apple and orange are not really royal.

>>3. Feature-age:

words man and woman doesn't connote much about age. Maybe men and woman implies that they're adults, but maybe neither necessarily young nor old. So maybe values close to 0. Whereas kings and queens are always almost always adults. And apple and orange might be more neutral with respect to age.

>>4. Feature-food:

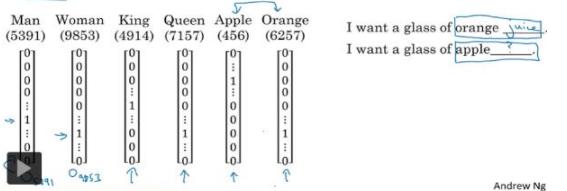
word man is not a food, woman is not a food, neither are kings and queens, but apples and oranges are foods.

These words can be many other features as well ranging from size, cost, Is this something that is a live, Is this an action, or is this a noun, or is this a verb, or is it something else? And so on. So you can imagine coming up with many features.

Word representation

$$\mathbf{V} = [\mathbf{a}, \mathbf{aaron}, \dots, \mathbf{zulu}, \mathbf{UNK}]$$

1-hot representation



网易云课堂

|V| = 10,000

Andrew Ng

1.2 featurized representation

>2. Feature note:

If take 300 different features, and it allows you to take this list of 300 numbers, that then becomes a 300 dimensional vector for representing the a word.

e.g.: take word man. use the notation \mathbf{e}_{5391} to denote a representation of its feature vector. while 5391 corresponding to word position in dictionary.

And similarly, for the other words.

>3. Summarize:

if use this representation-feature vector to represent the words orange and apple, then the representations for orange and apple are now quite similar: Some of the features will differ because of the color of an orange, the color an apple, the taste, but by a large, a lot of the features of apple and orange are actually the same, or take on very similar values. And so, this increases the odds of the learning algorithm that has figured out that orange juice is a thing, to also quickly figure out that apple juice is a thing.

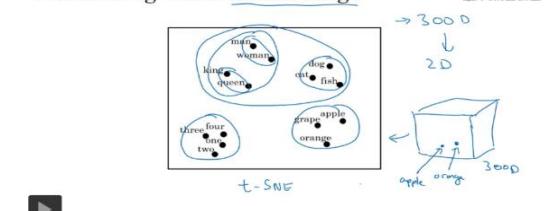
So feature vector representing allows algorithm to generalize better across different words.

Featurized representation: word embedding

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
... verb	\mathbf{e}_{5391}	\mathbf{e}_{9853}				

I want a glass of orange juice.
I want a glass of apple juice.
Andrew Ng

Visualizing word embeddings



网易云课堂

Andrew Ng

3. Feature visualization

Over the next few videos, we'll find a way to learn words embeddings: We just need to learn high dimensional feature vectors that gives a better representation than one-hot vectors for representing different words.

The features we'll end up learning, won't have an easy to interpret interpretation: will not be interpretation like that component one is gender, component two is royal, component three is age and so on. Exactly what they're representing will be a bit harder to figure out.

But nonetheless, the featurized representations we will learn, will allow an algorithm to quickly figure out that apple and orange are more similar than, king and orange or queen and orange.

3.1 Visualization-feature vector:

Once learn a feature vector (e.g: 300 dimensional), one of the popular things to do is also to take this 300 dimensional data and embed it in a 2D dimensional space to visualize them.

3.2 Visualizaiton method: one common algorithm for doing this is the t-SNE algorithm: (due to Laurens van der Maaten and Geoff Hinton).

3.3 Visualization example:

In right 2D visualization of words feature vectors: find that words like man and woman tend to get grouped together, king and queen tend to get grouped together, and these are the people which tends to get grouped together;

Those are animals who can get grouped together.

Fruits will tend to be close to each other. Numbers like one, two, three, four, will be close to each other.

And then, maybe the animate objects as whole will also tend to be grouped together.

This feature vector visualization gives a sense/intuition that, word embeddings algorithms like this can learn similar features for concepts that feel like they should be more related, as visualized by that concept that seem to you and me like they should be more similar, end up getting mapped to a more similar feature vectors.

3.4: word embedding naming:

Embedding: these representations will use these sort of featurized representations in maybe a 300 dimensional space, these are called embeddings.

The reason of calling them embeddings, is you can think of a 300 dimensional space, take every words like orange, and has a 300 dimensional feature vector so that word orange gets embedded to a point in this 300 dimensional space.

Visualization method:

Algorithms like t-SNE, map feature vector to a much lower dimensional space: can actually plot the 2D data and look at it.

that's what the term embedding comes from.

5.2-2 Using word embedding_NLP and word embedding

In the last video, you saw what it might mean to learn a featurized representations of different words. In this video, you see how we can take these representations and plug them into NLP applications.

1. Named entity recognition example:

1.1: Examples

Object-->Trying to detect people's names

> Sentence 1-training set: Given a test sentence like "Sally Johnson is an orange farmer"
Hopefully, trained algorithm figure out that Sally Johnson is a person's name, hence, the outputs is 1.
And one way to be sure that Sally Johnson has to be a person, rather than say the name of the corporation is that you know orange farmer is a person. (BRNN needed)

>Sentence 2-test set: New test sentence: "Robert Lin is an apple farmer"

if now use the featurized representations- the embedding vectors, then after having trained a model that uses word embeddings as the inputs, if has above new input, Knowing that orange and apple are very similar will make it easier for learning algorithm to generalize to figure out that Robert Lin is also a human, is also a person's name. (Word embedding needed as input)

>Sentence3-test set: has much less common words in sentence: " Robert Lin is a durian cultivator"

A durian is a rare type of fruit, if you have a small label training set for the named entity recognition task, you might not even have seen the word durian or seen the word cultivator in your training set.
But if you have learned a word embedding that tells you that durian is a fruit, so it's like an orange, and a cultivator, someone that cultivates is like a farmer, then you might still be generalize from having seen an orange farmer in your training set to knowing that a durian cultivator is also probably a person.

1.1 Word embedding able to do transfer learning is :

Word embedding enable algorithm to generalize words learned in training set to test set, due to word embedding

algorithm has a large training set of unlabeled text, examin very large text corpuses-get large words' feature.

(note: word embedding matrix is large enough, and trained by large unlabeled text-->get large words feature vector -which reflected words' similarity;

-->words/feature vectors used in new algorithm training set, could cover new algorithm test set (feature vector, as relation with training feature vector has been learned in word embedding matrix/feature vector per se)

if their feature vector is similar, then RNN algorithm could generalize from knowing one words role to another word of similar feature.

So one of the reasons that word embeddings will be able to do this (generalizing words learned in training set to test set) is the algorithms to learning word embeddings can examine very large text corpuses, maybe found off the Internet. So you can examine very large data sets, maybe a billion words, maybe even up to 100 billion words would be quite reasonable. So very large training sets of just unlabeled text.

And by examining tons of unlabeled text, which you can download more or less for free, you can figure out that orange and durian are similar. And farmer and cultivator are similar, and therefore, learn embeddings, that groups them together.
Now having discovered that orange and durian are both fruits by reading massive amounts of Internet text, what you can do is then take this word embedding and apply it to your named entity recognition task, for which you might have a much smaller training set.

---->

2. Transfer learning and word embedding

Process of carrying out transfer learning using word embeddings:

1. first is to learn word embeddings from a large text corpus (1-100B words): E
(or download pre-trained word embeddings online)

There are several word embeddings that you can find online under very permissive licenses.

2. Transfer embedding to new task, where you have a much smaller labeled training sets (like 100k words).

And use this like 300 dimensional embedding, to represent your words. One nice thing also about this is you can now use relatively lower dimensional feature vectors. So rather than using a 10,000 dimensional one-hot vector, you can now instead use maybe a 300 dimensional dense vector. Although the one-hot vector is fast and the 300 dimensional vector that you might learn for your embedding will be a dense vector

3. Optional: continue to fine tune the word embedding with new data

as you train your model on your new task with a smaller label data set, one thing you can optionally do is to continue to fine tune / adjust the word embeddings with the new data.

In practice, do this only if new task has a pretty big data set. If your label data set for new task is quite small, then usually, would not bother to continue to fine tune the word embeddings.

Transfer learning and word embeddings

1. Learn word embeddings from large text corpus. (1-100B words)

(Or download pre-trained embedding online.)

2. Transfer embedding to new task with smaller training set.
(say, 100k words)

$\rightarrow 10,000 \rightarrow 300$

3. Optional: Continue to finetune the word embeddings with new data.

3. Word embedding relation to face encoding

3.1: Face recognition:

Train a Siamese network architecture that would learn, say, a 128 dimensional representation for different faces. And then you can compare these encodings in order to figure out if these two pictures are of the same face.

The words encoding and embedding mean fairly similar things.

In the face recognition literature, people also use the term encoding to refer to these vectors, $f(x(i))$ and $f(x(j))$.

3.2: Difference of word embedding and face encoding:

One difference between the face recognition literature and what we do in word embeddings is:

> for face recognition, we want to train a neural network that can take as input any face picture, even a picture you've never seen before, and have a neural network compute an encoding for that new picture.

> Whereas for learning word embedding: we'll have a fixed vocabulary of, say, 10,000 words. And we'll learn a vector e_1 through, say, $e_{10,000}$ that just learns a fixed encoding or learns a fixed embedding for each of the words in our vocabulary

So that's one difference between the set of ideas you saw for face recognition versus what the algorithms we'll discuss in the next few videos.

Summary:

This item show how using word embeddings allows to implement this type of transfer learning. And how, by replacing the one-hot vectors we're using previously with the embedding vectors, you can allow your algorithms to generalize much better, or you can learn from much less label data.

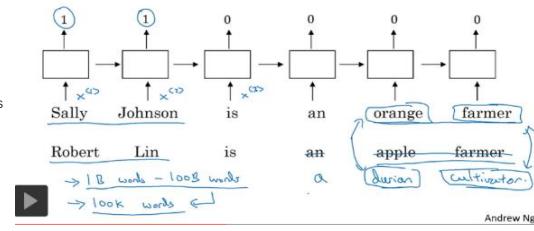
Next, I want to show you just a few more properties of these word embeddings. And then after that, we will talk about algorithms for actually learning these word embeddings. Let's go on to the next video, where you'll see how word embeddings can help with reasoning about analogies.

5.2.3 Properties of word embedding

One of the most fascinating properties of word embeddings is that they can also help with analogy reasoning. And while reasonable analogies may not be by itself the most important NLP application, they might also help convey a sense of what these word embeddings are doing, what these word embeddings can do.

Named entity recognition example

网易云课堂



1.1 Word embedding able to do transfer learning is :

--->

And so this allows to carry out transfer learning: where you take information you've learned from huge amounts of unlabeled text that you can suck down essentially for free off the Internet to figure out that orange, apple, and durian are fruits.

And then transfer that knowledge to a task, such as named entity recognition, for which you may have a relatively small labeled training set.

Note: for simplicity, in this example only draw a unidirectional RNN. But for named entity recognition task, should use BRNN.

(Note: Word embedding enable trained algorithm (with small training set) to transfer to test words: due to word embedding :)

trained word embedding algorithm get word embedding matrix: feature number x vocab size --> each word in vocab get its learned feature;

for algorithm trained on small training set, once learned training set - words' embedding vector as input, then could generalize to similar embedding vectors-->corresponding to test set (should have same distribution as training, dev set -> feature vectors in training set could almost cover feature vector in dev, test set.).

(Note: use smaller training set/feature vectors to represent whole data set (feature vector), as one training example/feature vector could represent ten/more examples(feature vector) in dev/test set. kind of like: pre-processing database , get their relationship by data feature vector-->one training set need to be more representative for the whole database, smaller training set need.)
(large training set needed as to cover almost all features of the whole database, as now database itself already learned relationships internally, no need so many training set).

2. Transfer learning and word embedding

Summarize

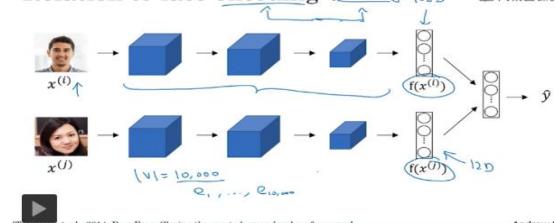
Word embeddings for transfer learning:

1. Tend to make the biggest difference when the task you're trying to carry out has a relatively smaller training set. So it has been useful for many NLP tasks like named entity recognition, text summarization, co-reference resolution, parsing.

2. Less useful for language modeling, machine translation, especially if you're accessing a language modeling or machine translation task for which you have a lot of data just dedicated to that task.

If trying to transfer from some task A to some task B, the process of transfer learning is just most useful when you happen to have a tons of data for A and a relatively smaller data set for B.
And so this rule is true for a lot of NLP tasks, but just less true for some language modeling and machine translation settings.

Relation to face encoding (embedding)



3. Word embedding relation to face encoding

Note:

the terms encoding and embedding are used somewhat interchangeably. So the difference of face encoding and word embedding is not represented by the difference in terminologies.

It's just a difference in how we need to use these algorithms in face recognition, where there's unlimited sea of pictures you could see in the future.

Versus natural language processing, where there might be just a fixed vocabulary, and everything else like that we'll just declare as an unknown word.

1. Analogies

e.g: man -> woman, king ->?

Hope algorithm could learn analogies automatically.

First calculate vector: $e_{\text{man}} - e_{\text{woman}}$, then find a word could get similar result vector by calculating with e_{king} : $e_{\text{king}} - e_{\text{woman}}$ ~ $e_{\text{king}} - e_{?}$

1.1 E.g: question, man is to woman as king is to_ what?

-> to have an algorithm figure this 'queen' out automatically?

vocabulary feature: only four feature: Gender, Royal, Age, Food.

Process to find word 'queen':

>1. using this four dimensional vector-feature vector to represent man, denote as e5391, and also woman, king, queen.

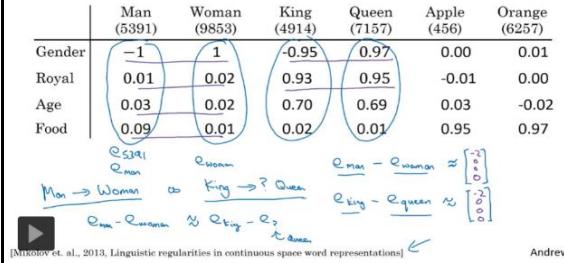
>2. One interesting property of these vectors is :

$$e_{\text{man}} - e_{\text{woman}} \approx [-2, 0, 0, 0]^T$$

$$e_{\text{king}} - e_{\text{queen}} \approx [-2, 0, 0, 0]^T$$

This is capturing is that the main difference between man and woman is the gender. And the main difference between king and queen, as represented by these vectors, is also the gender. Which is why the difference $e_{\text{man}} - e_{\text{woman}}$, and the difference $e_{\text{king}} - e_{\text{queen}}$, are about the same.

Analogies



[Mikolov et. al., 2013, Linguistic regularities in continuous space word representations]

网易云课堂

Andrew Ng

1. Analogies

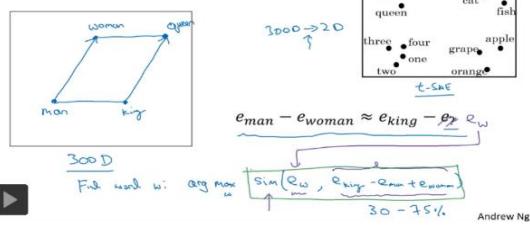
1.2: One way to carry out this analogy reasoning

If the algorithm is asked, man is to woman as king is to what?

Algorithm can compute vector difference of $e_{\text{man}} - e_{\text{woman}}$, and try to find a feature vector-a word, so it's difference with king feature vector is close to $e_{\text{man}} - e_{\text{woman}}$. And it turns out that when queen is the word plugged in here, then the left hand side is close to the right hand side.

So these ideas were first pointed out by Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. And it's been one of the most remarkable and surprisingly influential results about word embeddings.

Analogies using word vectors



2 Analogies using word vectors

in Feature space, vector: $e_{\text{man}} - e_{\text{wan}}$ is the vector difference, very similar to the vector difference; $e_{\text{king}} - e_{\text{queen}}$ --

>vector difference actually is a gender difference

Analogy method: word vector:

finding a new word w, to maximize the similarity : max sim (e_w , $e_{\text{king}} - e_{\text{man}} + e_{\text{woman}}$)

have some appropriate similarity function for measuring how similar is some word w to this quantity of the right.

Note:

1. this method - using word vectors works , by maximizing the similarity, can actually get the exact right answer.
2. Accuracy in papers like 30%-75% , count an analogy attempt as correct only if it guessed the exact word right.

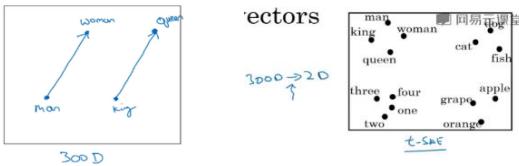
2.3 feature vector similarity in visualized space

Algorithms like t-SAE could visualize words from high dimension feature space to lower dimension space.

SAE : it takes original feature vector (like 300-D data), and maps it in a very non-linear way to a 2D space.

The mapping that t-SAE learns is a very complicated and very non-linear mapping. So after the t-SAE mapping, should not expect these types of parallelogram -words' feature vector difference relationships, like the one we saw on the left, to hold true.

And it's really in this original 300 dimensional space that you can more reliably count on these types of parallelogram relationships in analogy pairs to hold true. And it may hold true after a mapping through t-SAE, but in most cases, because of t-SAE's non-linear mapping, you should not count on that. And many of the parallelogram analogy relationships (平行四边形类似关系) will be broken by t-SAE.



4. similarity function that is most commonly used

The most commonly used similarity function is called cosine similarity.

In cosine similarity, define the similarity between two vector as basically the inner product between u and v.

cosine similarity: cosine of the angle between two vectors: $\text{Sim}(u, v) = u \cdot v / (\|u\|_2 \times \|v\|_2)$

if u, v are very similar, inner product $u \cdot v$ tend to be large: angle = 0 -> consie=1, angle=90, consie=0, angle =180, cosine=-1

This is called cosine similarity because this is actually the cosine of the angle between the two vectors, u and v.

Note:

1. Cosine similarity works quite well for these analogy reasoning tasks.

2. Can also use square distance or Euclidean distance, $-\|u - v\|^2$.

Technically, this would be a measure of dissimilarity rather than a measure of similarity. So we need to take the negative of this, and this will work okay as well. Yet cosine similarity is being used a bit more often.

The main difference between these similarity functions is how it normalizes the lengths of the vectors u and v.

Summary:

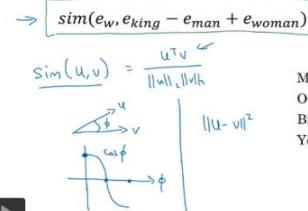
One of the remarkable results about word embeddings is the generality of analogy relationships they can learn: e.g , it can learn that man is to woman as boy is to girl, because the vector difference between man and woman, similar to king and queen and boy and girl, is primarily just the gender.

And all of these things can be learned just by running a word embedding learning algorithm on the large text corpus. It can spot all of these patterns by itself, just by running from very large bodies of text.

So in this video, you saw how word embeddings can be used for analogy reasoning. And while you might not be trying to build an analogy reasoning system yourself as an application, this I hope conveys some intuition about the types of feature-like representations that these representations can learn. And you also saw how cosine similarity can be a way to measure the similarity between two different word embeddings. Now, we talked a lot about properties of these embeddings and how you can use them. Next, let's talk about how you'd actually learn these word embeddings, let's go on to the next video

(note: word embedding matrix: -->transfer learning (: new algorithm trained with less training data) + analogy reasoning (similarity function of words embedding vecots)

Cosine similarity



Man:Woman as Boy:Girl
Ottawa:Canada as Nairobi:Kenya
Big:Bigger as Tall:Taller
Yen:Japan as Ruble:Russia

5.2-4 embedding matrix

Formalize the problem of learning a good word embedding. When you implement an algorithm to learn a word embedding, what you end up learning is an embedding matrix

1. Embedding matrix
e.g: using 10,000-word vocabulary, each word have 300 feature.
Object: s learn embedding matrix E, 300 dimensional x 10,000 dimensional matrix.

1.1 Embedding matrix :E (feature number x vocabulary size)
Columns: different embeddings for the each word in vocabulary.
Rows: feature for words

1.2 One-hot vector: O_j
vector element is 0/1, corresponding to the position in vocabulary, 0 for other positions.
e.g.: Orange was word number 6257 in vocabulary of 10,000 words. So use notation O₆₂₅₇ as its one-hot vector, with zeros everywhere and a one in position 6257. And so, this will be a 10,000-dimensional vector with a one in just one position

1.3 for word j embedding vector : e_j , have equation: e_j = E · O_j
So, the embedding matrix E times one-hot vector here winds up selecting out this 300-dimensional column corresponding to the word j.

2. Learn an embedding matrix E
Idea: Initialize E randomly and use gradient descent to learn all the parameters of this feature number x vocabulary size (300 x 10,000) dimensional matrix.
and E · O_j this one-hot vector gives the embedding vector.

Note: equation e_j = E · O_j
Writing the equation, as it is convenient to write this type of notation.
But when implementing this, it is not efficient to actually implement this as a matrix vector multiplication, because the one-hot vectors, now this is a relatively high dimensional vector and most of these elements are zero.
In practice use a specialized function to just look up a column of the Matrix E rather than do this with the matrix multiplication.

Embedding matrix
Diagram illustrating the computation of an embedding vector. An input word 'orange' is mapped to a one-hot vector O₆₂₅₇, which is then multiplied by the embedding matrix E to produce the embedding vector e₆₂₅₇. The matrix E has 300 columns, corresponding to the 300 features. The row index for 'orange' is 6257. The resulting vector e₆₂₅₇ is also labeled as 'embedding for word j'. A summary note says: 'practice, use specialized function to look up an embedding.' Andrew Ng

5.2-5 Learning word embedding
In this item start to learn some concrete algorithms for learning word embeddings.
In the history of deep learning as applied to learning word embeddings, people actually started off with relatively complex algorithms. And then over time, researchers discovered they can use simpler and simpler and simpler algorithms and still get very good results especially for a large dataset. But what happened is, some of the algorithms that are most popular today, they are so simple that if I present them first, it might seem almost a little bit magical, how can something this simple work? So, here to start off with some of the slightly more complex algorithms it's actually easier to develop intuition about why they should work, and then we'll move on to simplify these algorithms and show you some of the simple algorithms that also give very good results.

1. Neural language model
e.g using neural network to build a language model, given previous few words, to predict the next word.
Input/training set: "I want a glass of orange."
expected output: during training, want neural network to predict the next word in the sequence.

Note: each word of the input already have the index in the vocabulary .

Building a **neural language model** is a reasonable way to learn a set of embeddings
(note: why not sequence language model? input word short, neural network more fast?
neural language model is just one time step of sequence language mode:
in time t, input x_{t-1} = correct y_{t-1}, a_{t-1}->-previous words activation:
 $a^{t-1} = g_a(W_{a^{t-1}} \cdot x^{t-1} + b_a)$
 $y^t = softmax(W_y \cdot a^{t-1} + b_y)$
-->loss function: max output probability of right label y_t

while transfer this step to neural network:
a<t> = g_a (W₁*[e₁, e₂...e_n] + b₁)
y^t<t> = softmax (W₂* a<t> + b₂)
)

Neural language model
Diagram illustrating a neural probabilistic language model. It shows the input sequence "I want a glass of orange" being processed through an embedding matrix E to produce word vectors e₁ through e₆₂₅₇. These vectors are then passed through a softmax layer with weights W₀₁ and bias b₀₁ to produce the probability distribution for the next word, which is "juice". The diagram also shows the hidden state vector h₀ and the final output word "apple juice". Andrew

1.1 Procedure : build a neural network to predict the next word in the sequence:
>1. Construct one-hot vector corresponding to each input word i. O_j
one-hot vector is vocabulary size dimensional vector.

>2. Get embedding vector for each input word: e_j: O_j----->E ----->e_j
 $e_j = E \cdot O_j$

>3. Feed all of these embedding vectors into a neural network-have its own parameter: W1, b1 :a<t> = W1 · e<t> + b1
each word embedding vectors will be stacked together and feed to softmax.
Input to neural network, dimension is : word embedding-feature number x input word number
e.g: 6 words, and feature number is 300, then neural network input dimension is 300 x 6-->unrolled (1800,1)
4. feeds to a softmax, which has it's own parameters as well: W2, b2; y_t<t> = g_y (W₂ · a<t> + b₂)
soft max classifies 1000 possible outputs in the vocabulary for the final word trying to predict.
Output vocabulary size vector: possibility of each word in vocabulary being the next word in input sentence.
-->so if in the training slide we saw the word juice then, the target for the softmax in training would be that it should predict the same word juice .

(note: could use RNN: many to one architecture??:
e<1> --> a<1>: g1(w_{a<1>} · a<1>) + b_{a<1>} --> + e<2> --> a<2> -->a<t> --> y^t = g2(w_y · a<t> + b_y)
parameter: only wa, wy, by, less than above; no: wa: dimension is a_activation dim x (a+n_feature dim),
while above w1: word number x n_feature dim, less parameters
but this still have input dimension change problem?-->**W1 dimension need to fix input size**

Neural language model
Diagram illustrating a neural probabilistic language model. It shows the input sequence "I want a glass of orange" being processed through an embedding matrix E to produce word vectors e₁ through e₆₂₅₇. These vectors are then passed through a softmax layer with weights W₀₁ and bias b₀₁ to produce the probability distribution for the next word, which is "juice". The diagram also shows the hidden state vector h₀ and the final output word "apple juice". Andrew

1.2 Note:
> 1. Historical window size
word embeddings feed to neural network: More commonly have a fixed historical window-fixed window size, is hyperparameter of algorithm.
e.g.: might decide always want to predict the next word given the previous four words (note: four here is a hyperparameter of the algorithm.)
So this is how you adjust to either very long or very short sentences where decide to always just look at the previous four words, so neural network will input a 1,200 dimensional feature vector, go into this hidden neural layer, then have a softmax and try to predict the output.

Using a fixed history, means can deal with even arbitrarily long sentences because the input sizes are always fixed.

>2. parameters in this language model:
>>1. Embedding matrix E: use same matrix E for all the words in embedding vector computation: e_j = E · O_j
don't have different matrices for different positions in the proceedings four words, is the same matrix E.
>>2. Weight W1, b1; W2, b2
These parameters of the algorithm can use backprop to perform gradient descent to maximize the likelihood of your training set : e.g repeatedly predict given four words in a sequence, what is the next word in vocabulary.

1.3 Summarize:
This algorithm will learn pretty decent word embeddings.
Reason: if algorithm remember orange juice, apple juice example, it is in the algorithm's incentive to learn pretty similar word embeddings for orange and apple, because doing so allows it to fit the training set better:
as it's going to see orange juice sometimes, or see apple juice sometimes, and so the algorithm will find that it fits the training set fast if apples, oranges, and grapes, and pears, and so on and maybe also durians , with similar feature vectors.
So, this is one of the earlier and pretty successful algorithms for learning word embeddings- this matrix E.

(note: this is one time step of language model: stack this model here in time step -->get the language model:
language model: predict sentence probability: p(y<1>..y<Ty>) = p(y<1>) * p(y<2>|y<1>) * p(y<3>|y<1>, y<2>)...
except p(y<1>): generated by x<1> = 0, a<1>=0, others time step input is y<t-1>, a<t-1>;
here in word embedding training, a<t-1> is all previous word e₁..e<t-1>, not the activation of previous words)

2. Other context/target pairs

Generalize above algorithm and derive even simpler algorithms.

e.g. input sentence in training set : " I want a glass of orange juice to go along with my cereal".
Object of the algorithm was to predict some target word 'juice' , given some context -like the last four words.

2.1 Context to use for training algorithm

Researchers have experimented with many different types of context and get this conclusion:

>1. If the goal is to build a language model:

Then it is natural for the context to be a few words right before the target word.

>2. If the goal is not to learn the language model per se,

Then can choose other contexts:

>>>1. 4 words on left and right

Means we're posing a learning problem where the algorithm is given four words on the left, and four words on the right, to go along with, and this has to predict the word in the middle.

And posing a learning problem like this where you have the embeddings of the left four words and the right four words feed into a neural network, similar to try to predict the word in the middle, try to put it target word in the middle.

This can also be used to learn word embeddings.

3. Summary:

In this video show how the language modeling problem which causes the pose of machines learning problem: where you input the context like the last four words and predicts some target words. How posing that problem allows you to learn input word embedding.

In the next video, you'll see how using even simpler context and even simpler learning algorithms to map from context to target word, can also allow you to learn a good word embedding.

2. Other context/target pairs

2.1 Context to use for training algorithm

>>>2. Last 1 word:

This will be different learning problem: given one word, and predict the next word.
Can construct a neural network that just feed the embedding of the one previous word to a neural network as you try to predict the next word.

>>>3. Nearby one word-works surprisingly well

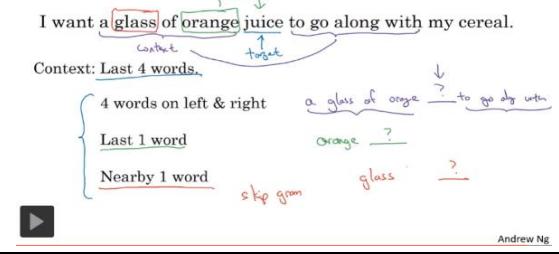
e.g. take the word 'glass', and then there's another words somewhere close to glass, what do you think that word is?

This is the idea of a Skip-Gram model, a simpler algorithm where the context is just one word rather than four words, but this works remarkably well.

Summary:
If you really want to build a language model, it's natural to use the last few words as a context.
But if your main goal is really to learn a word embedding, then can use all of these other contexts and they will result in very meaningful word embeddings as well.

(note: word embedding transfer learning is less useful for language model , which has enough large data to train algorithm, -->learn word embedding by simpler method would be better??)

Other context/target pairs



网易云课堂

Andrew Ng

5.2-6 Word2Vec_NLP and word embeddings

In the last video, showed how by learning a neural language model to get good word embeddings.

In this item, will see the Word2Vec algorithm which is simple and comfortably more efficient way to learn this types of embeddings.

1. Skip-grams

E.g:

Given one sentence of training set : "I want a glass of orange juice to go along with my cereal"

Object: In the skip-gram model, come up with a few context to target pairs to create our supervised learning problem.

1.1: Idea:

>Most of the ideas here are due to Tomas Mikolov, Kai Chen, Greg Corrado, and Jeff Dean.

1.2: Set a supervised learning problem

Process:

>1. Context word: Rather than having the context be always the last four words or the last end words immediately before the target word, instead randomly pick a word to be the context word.

e.g chose the word orange as context word

>2. Target word: Randomly pick another word within some window of chosen context word as target word.

e.g. within window +/- 5 words or +/- 10 words of the context word , randomly choose a word to be target word.

So maybe just by chance you might pick juice to be a target word, Or might choose glass so have another pair , or maybe choose the word 'my' as the target.

And so we'll set up a supervised learning problem where given the context word, you're asked to predict what is a randomly chosen word within like, +/- 10 or +/- 5 word window, of that input context word.

Note:

This is not a very easy learning problem, because within +/- 10 words of the word orange, it could be a lot of different words. But a goal that's setting up this supervised learning problem, isn't to do well on the supervised learning problem per se, it is that we want to use this learning problem to learn good word embeddings.

2. Model

e.g. vocab has 10,000 words. S

In vocabulary, orange is word O_6257, and the word juice is the word O_4834

Object: the basic supervised learning problem we're going to solve is that we want to learn the mapping from some Context c, such as the word orange to some target, which might be the word juice or the word glass or the word my.

(note: labeled y: only label one word in the specify window of context word as right(1), other word in vocab as 0.

but in training set, there are many other words in the specify window of context word, if algorithm learned to output possibility for those word close to 0, should be a problem??

input x: context word
labelled output y: target word t

Object: learn to map x to that open y.

2.1 Model process

>1. Construct input context word: get embedding vector for the input context word

input word -->one-hot vector O_c ----->embedding matrix E ----->embedding vector e_c = E.O_c

>2. Feed context word pembedding vector to neural network.

(note: similar to use neural language model as above, only directly go into softmax, not middle activation)

2.2 Mode-detail. softmax model:

Probability of different target words given the input context word :

$$p(t|c) = \exp [(\theta_t)^T * e_c] / \sum (\exp [(\theta_j)^T * e_c]) \text{ over all word in vocab.}$$

2.3 . loss function for softmax

$$\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{10,000} y_i \log \hat{y}_i$$

target word:

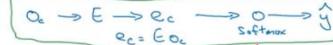
>labeled y: represented by one-hot vector

>predicted \hat{y}^n , vector of vocab size, probability of being the target for each word in vocab.

$$l(y^n, \hat{y}) = - \sum y_i \log (\hat{y}^n_i) \text{ over all word in vocab.}$$

e.g.:

if the target word is juice O_4834, then its one-hot vector 4834-th element equal to 1 and the rest is 0. And similarly \hat{y}^n will be a 10,000 dimensional vector output by the softmax unit with probabilities for all 10,000 possible targets words.-->loss will be - log y_{4834}^n . (y_{4834}^n word juice's probability for being the target word)



2.4 Summarize:

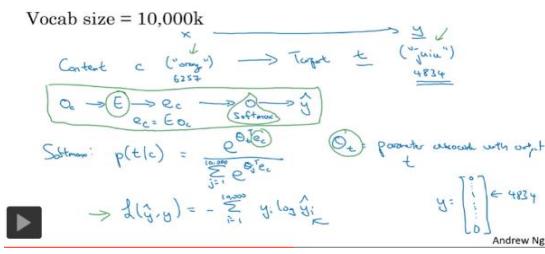
Above is the overall little model, - neural network with basically looking up the embedding and then just a soft max unit.

The matrix E will have a lot of parameters corresponding to all of these embedding vectors,e_c and the softmax unit also has parameters θ_i ($i = 1, \dots, 10,000$ -vocab. size)

by optimizing this loss function with respect to all of these parameters,will actually get a pretty good set of embedding vectors.

This is called the skip-gram model: because it is taking as input one word like orange and then trying to predict some words skipping a few words from the left or the right side. To predict what comes little bit before/ little bit after the context words.

Model



3. Problems with softmax classification

$$p(t|c) = \frac{e^{\theta_c^T e_t}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_t}}$$

There are a couple problems with using this algorithm, and the primary problem is computational speed

3.1 Softmax classification computation cost:

Every time/iteration to evaluate this probability, need to carry out a sum over all vocab words : sum $((\theta_j)^T \cdot e_c)$ over all word in vocab. (θ_j different in every iteration)

when vocab size is very big like 100,000 or a 1,000,000, it gets really slow to sum up over this denominator every single time/iteration/step. In fact 10,000 is actually already quite slow, and it makes even harder to scale to larger vocabularies.

3.2 hierarchical softmax classifier

One solutions to this is to use a hierarchical softmax classifier.



hierarchical softmax classifier: instead of trying to categorize something into all vocabulary words like 10,000 categories on one go but:

>1. Have first binary classifier, tells if the target word is in the first half vocab, or in the second half of vocab.
If this binary cost tells target word in the first half vocab like first 5,000 words, then use a second binary classifier.

>2. use a second binary classifier, to tell if the target word in the first half or second half of 'first half vocab'

continue do this: using binary classifier to localize target word in a smaller new 'vocab size', until eventually you get down to classify exactly what word it is.

4. Sampling the context C

Once sampled the context C, the target t can be sampled within some word window of the context C.

Choosing context C: method:

>1. Sample uniformly, at random, from training corpus

Not recommended , as training site will be dominated by these extremely frequently words, because then you spend almost all the effort updating e_c , for those frequently occurring words.

e.g: there are some words like 'the, of, a, and, to' and so on that appear extremely frequently, if sampling uniformly, at random, then context to target mapping pairs just get these types of words extremely frequently, whereas there are other words like orange, apple, and also durian that don't appear that often.

But we need to make sure to spend some time updating the embedding, even for these less common words like e_durian.

>2. heuristics to balance common and less common word

In practice the distribution of words p(c) isn't taken just entirely uniformly at random for the training set corpus, but instead there are different heuristics that you could use in order to balance out something from the common words together with the less common words.

HOW TO IMPLEMENT IN CODING?

How to generate training set?

5.2-7 Negative sampling_NLP and word embedding

Last video showed how the Skip-Gram model allows to construct a supervised learning task: mapping from context to target and this allows you to learn a useful word embedding.

But the downside of that was the Softmax objective was slow to compute. (should be a common problem for algorithm using softmax with large output dimension)

In this video, will show a modified learning problem - negative sampling: do something similar to the Skip-Gram model but with a much more efficient learning algorithm.

1. Defining a new learning problem

Most of the ideas presented in this video are due to Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeff Dean.

1.1 negative sampling algorithm object: create a new supervised learning problem: given a pair of words like orange and juice, to predict if this is a context-target pair.

1.2: example:
e.g: training sentence: 'I want a glass of orange juice to go along with my cereal'

Input: orange - juice was a positive example, target (output) $y = 1$
Input: orange - king is a negative example, target (output) $y = 0$

1.3: Generating training set:

1.3.1 Process:

>1. generated a positive input/ pair: sample a context and a target word.

e.g: orange - juice is associate that with a label of 1

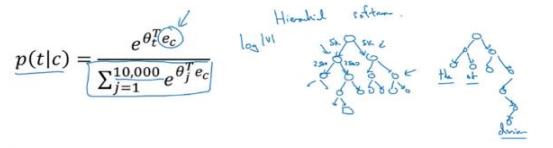
Positive example generation method: generated exactly as skip model: in a training sentence, Sample a context word, look around a window of say, plus-minus ten words and pick a target word.

>2. generate k negative examples:

Take the same context word and then just pick a word at random from the dictionary.
Under the assumption that if we pick a random word from vocab, it probably won't be associated with the word from training sentence: so

e.g: chose the word king at random and label pair orange-king as 0. and then take orange and pick another random word from the dictionary, and label this pair as 0.

Problems with softmax classification



How to sample the context c?

\rightarrow the, of, a, and, to, ...

\rightarrow orange, apple, durian

$c \rightarrow t$

$p(c)$

Andrew Ng

3. Probes with softmax classification

Note: Hierarchical softmax:

1. Having a tree of classifier like right, means that each of the interior nodes of the tree can be just a binding classifier-like logistic classifier.
So don't need to sum over all vocab size in order to make a single classification.

2. Computational classifying tree like this, scales like log of the vocab size rather than linear in vocab size.
That is why it is called a hierarchical softmax classifier.



3. In practice, the hierarchical softmax classifier doesn't use a perfectly balanced tree or this perfectly symmetric tree, with equal numbers of words on the left and right sides of each branch.
In practice, the hierarchical softmax classifier can be developed: common words tend to be on top, whereas the less common words like durian can be buried much deeper in the tree:
Because we see the more common words more often, and so in this way need only a few traversals to get to common words like 'the' and 'of'. Whereas less seen frequent words like durian is okay to be buried deep in the tree as see it less often-do not go that deep frequently.

--> So there are various heuristics for building the tree how you used to build the hierarchical softmax classifier.

In the next video, will talk about a different method- nectar sampling, which I think is even simpler. And also works really well for speeding up the softmax classifier and the problem of needing the sum over the entire vocab size in the denominator

5. Summary:

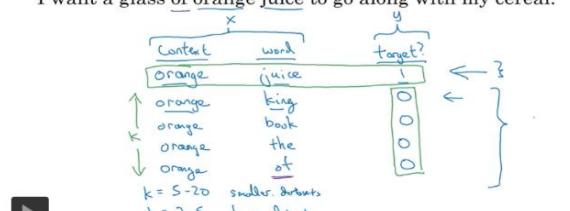
Above is for the Word2Vec skip-gram model. The original paper actually had two versions of this Word2Vec model, the skip gram was one, and the other one is called the CBow, the continuous backwards model, which takes the surrounding contexts from middle word, and uses the surrounding words to try to predict the middle word, and that algorithm also works, it has some advantages and disadvantages.

But the key problem with this algorithm with the skip-gram model as presented so far is that the softmax step is very expensive to calculate because needing to sum over entire vocabulary size into the denominator of the soft max.

In the next video I show you an algorithm that modifies the training objective that makes it run much more efficiently therefore lets you apply this in a much bigger fitting set as well and therefore learn much better word embeddings.

Defining a new learning problem

I want a glass of orange juice to go along with my cereal.



Andr

1.3.2 Note:

for all pairs, whose target word sampling from vocab instead of training sentence that context word in, labeled as 0, even if the target word appears next to context word in training sentence as well.
e.g: sampled target word 'of' from vocab, label orange-of as 0 also even if 'of' next to orange in training sentence.

1.3.3 Summarize:

for generating training data-positive and negative pairs:
Pick a context word and then pick a target word in training sentence - the first row of this table in right pic., this gives us a positive example. So context+target, and then give that a label of 1.

For some number of times say, k times: take the same context word and then pick random words from the dictionary, king, book, the, of, whatever comes out at random from the dictionary and label all those 0, and those will be our negative examples. and it's okay if just by chance, one of those words we picked at random from the dictionary happens to appear in the window, in a plus-minus ten word window say, next to the context word, orange. (also label as 0)

(note: could use logistic model??, architecture too simple to learn word embedding?)

train in set (c, t), labeled 0/1?-->input is a pair (c,t)-->sigmoid output 0/1

2. Mode:

Supervised learning model for learning a mapping from x to y.

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

Input/Training set: one example: positive and negative pair, + target labeled 1 or 0

input context word embedding e_c , labeled output $y = [0, 0, \dots, 0, 1]$ -vocab size/K+1 size.

output: predict for each pair, target is 0 or 1. (vocab size vector: for each word is it positive target 1 or negative target 0)

Denotion:

C: context word,

t: possible target word,

target: denote 0, 1, this is for a context target pair.

$$P(y=1 | c, t) = \sigma(\theta_c^T e_c)$$

2.1 Mode:

Define a logistic regression model: chance of $y = 1$, given the input c, t pair, $p(y=1|c, t) = \text{sigmoid}(\theta T \cdot e_c)$

$y^* = \text{sigmoid}(w^T e_c + b)$:

Input size: e_c : feature size 1

Output size: y^* : vocab size $\times 1$

Weight: $\rightarrow w$: dim: $v \times f = [w_1; w_2; \dots; w_v]$

w_t : parameter of each vocal word-> θ_t

2.2 Model parameter: θ_t, e_c

Vector θ_t for each possible target word t.

Embedding vector e_c for each possible context word.

2.3 Model illustration:

For every positive examples, have k negative examples with which to train this logistic regression-like model.

>1. Generate context word embedding vector

context word C as input--->one-hot vector O_c ----->E----->context word embedding vector: e_c

2. Feed context word embedding vector to softmax: $y = \text{sigmoid}(\theta T \cdot e_c)$

Output y^* : vector (vocab. size $\times 1$): predicting for each word t in vocab, the chance of $y=1 | c, t$

-->predict for vocab size pairs (c, t): t take each word from vocab, get a pair (c, t)-c is same, and predict their target 0/1.

Now it is vocab size like 10,000 logistic regression classification problems: each classifier (θ_t) corresponding to predict if the target word is this specific word t like juice in vocab.

So think of this as having vocab size (like 10,000) binary logistic regression classifiers.

2.4 . parameter updating:

But instead of training all 10,000 of them on every iteration, we're only going to train $k+1$ (1 positive and k negative=4 here) of them.

Train the one responding to the actual target word and then train k number randomly chosen negative examples.

Model: input context word --->one-hot O_c --->E--->embedding vector e_c ----->sigmoid on vocab : vocab size classifier: -->get vocab size output vector .

for one word t, in vocab:

>1. forward prop:

$Z_t = \theta T \cdot e_c$

$y^t = a_t \cdot Z_t = \text{sigmoid}(\theta T \cdot e_c)$

word t prediction loss: $L_t = L(y^t, y_t) = y_t \cdot \log(y^t) + (1-y_t) \cdot \log(1-y^t)$

>2. backprop:

$dZ_t = -y_t^* \cdot (1-y_t) + y^t \cdot (1-y_t) = y^t \cdot t - y_t$

$d\theta_t = dZ_t \cdot e_c$

$d e_c = dZ_t \cdot \theta_t$

>3. update parameter:

$\theta_t += -a \cdot d\theta_t$

$e_c += -a \cdot d e_c$

as only $k+1$ labeled y_t - if it is not the word of negative and positive pairs, no label y,-->no loss L could compute--> then could only do backprop for these labeled $k+1$ words.

3.Selecting the negative examples

So after having chosen the context word orange, how do you sample these words to generate the negative examples

>1. Sample according to the empirical frequency of words in your corpus:

So just sample it according to how often different words appears. But the problem with that is that will end up with a very high representation of words like 'the', 'of', and, and so on.

>2. The other extreme: use 1/(vocab size), sample the negative examples uniformly a $P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$

but this is also very non-representative of the distribution of English words

>3. Authors, Mikolov et al, reported that empirically, this heuristic value work best:

A little bit in between the two extremes of sampling from the empirical frequencies, meaning from whatever's the observed distribution in English text to the uniform distribution.

Sampled proportional to their frequency of a word to the power of 3/4:

e.g: $f(w)$ is the observed frequency of a particular word in the English language or in your training set corpus, then by taking it to the power of 3/4, this is somewhere in-between the extreme of taking uniform distribution. And the other extreme of just taking whatever was the observed distribution in your training set.

And so I'm not sure this is very theoretically justified, but multiple researchers are now using this heuristic, and it seems to work decently well.

1.4: Create a supervised learning problem

>1. learning problem:

>> training set:
Inputs x: these pair of words, positive and negative pairs: input context word embedding e_c
labeled output--vocab size vector: only labeled 1 for target position, 0 for negative target word position. other position do not care: $[0, 0, \dots, 1, 0, 0]$; position in vocab.

or just make output $k+1$ dimension vecot: 1 for positive target, 0 for k negative target word .
0 for negative target word.

>2. Idea & object:

So the problem is really, given a pair of words like orange and juice:
do you think they appear together/ Do you think these two words is gotten by sampling close to each other . Or do you think one word is gotten from the training text and one word chosen at random from the dictionary?

algorithm object is really to try to distinguish, for a given pair of words, from which of the two distributions they come from: distribution one-label 1 distribution: close to each other; distribution two-label 0 distribution: not close at all-one from training sentence, one from vocab.

So this is how you generate the training set.

Note: choose k:-number of negative pairs for the same context word: Large values of k for smaller data sets, why?-->make training set larger: as Context word C number is small, increase K-->increase training pairs.

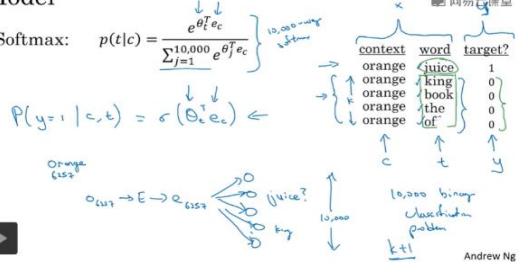
Mikolov et al, recommend :

for smaller data set: maybe k is 5 to 20,

If have a very large data set: then chose k to be smaller. k is like 2 - 5 for larger data sets. In this example, just use k = 4.

Model

Softmax: $p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$



2.5 Summarize:

So instead of having one giant 10,000 way Softmax, which is very expensive to compute, we've instead turned it into vocab size 10,000 binary classification problems, each of which is quite cheap to compute. And on every iteration, we're only going to train five of them or more generally, $k+1$ of them, of k negative examples and one positive examples.

(word2v-skim: -->one iteration will update all vocab word parameter θ_i & word embedding e_c)

And this is why the computation cost of this algorithm is much lower because you're updating $k+1$ units, $k+1$ binary classification problems. Which is relatively cheap to do on every iteration rather than updating a 10,000 way Softmax classifier.

This technique is called negative sampling:

Because what you're doing is, you have a positive example, the orange and then juice. And then you will go and deliberately generate a bunch of negative examples, negative samplings, hence, the name negative sampling, with which to train four more.

And on every iteration, choose four different random negative words with which to train your algorithm on.

(one training example every iteration??)

Selecting negative examples

context	word	target?
orange	juice	1
orange	king	0
orange	book	0
orange	the	0
orange	of	0

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$$

the , of, and...

$$\frac{1}{100}$$

4. Summarize

Learn word vectors in a Softmax classifier is very computationally expensive, by changing that to a bunch of binary classification problems, you can very efficiently learn a pretty good words vectors.

And as is the case in other areas of deep learning as well, there are open source implementations. And there are also pre-trained word vectors that others have trained and released online under permissive licenses. And so if you want to get going quickly on a NLP problem, it'd be reasonable to download someone else's word vectors and use that as a starting point.

So that's it for the **Skip-Gram model (negative sampling and word 2vect)**. In the next video, I want to share with you yet another version of a word embedding learning algorithm that is maybe even simpler than what you've seen so far.

(note:

word2 vect: context word -->target word:softmax (vocab size)

negative sampling: context word -->target positive and negative word-sigmoid (vocabsize, yet only train k+1 embeddings)-why need combine together 1+k negative? could not use context->target: sigmoid?-update only one word embedding every iteration-too slow so choose k negative? -->or could choose k' positive + k negative, for faster training??)

5.2-8 Glove word vectors_NLP and word embedding

You learn about several algorithms for computing words embeddings. Another algorithm that has some momentum in the NLP community is the GloVe algorithm. This is not used as much as the Word2Vec or the skip-gram models, but it has some enthusiasts, may because in part of its simplicity.

1. GloVe: (global vectors for word representation)

The GloVe algorithm was created by Jeffrey Pennington, Richard Socher, and Chris Manning.

1.1 Glove Context and target words:

Previously, we were sampling pairs of words, context and target words, by picking two words that appear in close proximity to each other in text corpus.

Glove algorithm: starts off just by making that explicit:

X_{ij} : is the number of times that a word i appears in the context of j.

word i play the role of word t, and j play the role of word c. so you can think of X_{ij} as being X_{tc} .

for each word i,j in vocab size. Go through your training corpus and just count up how many times does a word i appear in the context of a different word j.

(note: training set size: vocab. size x vocab size

if ij symmetric-->training set size = vocab size x vocab size / 2

GloVe (global vectors for word representation)

I want a glass of orange juice to go along with my cereal.

$$\begin{aligned} & c, t \\ & X_{ij} = \# \text{times } i \text{ appears in context of } j. \\ & X_{ij} = X_{ji} \leftarrow \end{aligned}$$

[Pennington et. al., 2014. GloVe: Global vectors for word representation]

Andrew Ng

2. Model:

2.1 Labeled training set:

Input: pair word-j-i; (note: training set size: vocab. size x vocab size; if ij symmetric-->training set size = vocab size x vocab size / 2)

labeled output: x_{ij} : How related are words i and j as measured by how often they occur with each other in training data.

2.2 model output:

>>> 1. $(\theta_i)^T e_j$: tell how related are those two words.

Think of i and j as playing the role of t and c--> $(\theta_i)^T e_j$, $(\theta_i)^T e_j$, is to tell you how related are those two words.

2.3 GloVe model object/ loss function:

Minimize the difference between $\theta_i^T e_j$ and $\log X_{ij}$

cost for a singal example: word j: $I(x_{ij})$: sum $(\theta_i)^T e_j - \log X_{ij}$ over all i in vocab;

cost for all training example : word j go through vocab size: sum $I(x_{ij})$ over all word j in vocab

2.4 Model parameter:

>>> 1. θ_i, e_j ($i=1, \dots, \text{vocab size}$)

Using gradient descent to minimize the sum difference of $(\theta_i)^T e_j - \log X_{ij}$ over i from 1 to vocab size, sum over j from 1 to vocab size;

to learn vectors, so that their end product $(\theta_i)^T e_j$ is a good predictor for how often the two words occur together, which represented here by x_{ij} .

(note: paramter θ_i, e_j play same role if X_{ij} is symmetric)

2.5 Note:

parameter θ_i, e_j ($i=1, \dots, \text{vocab size}$), their role are completely symmetric.

θ_i, e_j are symmetric in that, at the math, they play pretty much the same role and you could reverse them or sort them around, and they actually end up with the same optimization objective.

One way to train the algorithm is: initialize θ_i and e_j both uniformly at random, run gradient descent to minimize its objective, and then when you're done for every word, to then take the average.

e.g.: For a given words w, we have :

$e_w^{\text{final}} = e_w + \theta_w / 2$

e.w: word w embedding that was trained through this gradient descent procedure

$\theta_w: \theta_w$ trained through this gradient descent procedure

because θ and e in this particular formulation play symmetric roles unlike the earlier models we saw in the previous videos, where θ and e actually play different roles and couldn't just be averaged like that.

(note: skip-gram: negative sample * word2v: de_c, d8_w not same-->not play symmetric role)

That's it for the GloVe algorithm. one confusing part of this algorithm is, it seems almost too simple. How could it be that just minimizing a square cost function like this allows you to learn meaningful word embeddings? But it turns out that this works. And the way that the inventors end up with this algorithm was, they were building on the history of much more complicated algorithms like the newer language model, and then later, there came the Word2Vec skip-gram model, and then this came later. And we really hope to simplify all of the earlier algorithms.

3. A note on the featurization view of word embeddings

We started off with right talbe featurization view as the motivation for learning word vectors:

maybe the first component of the embedding vector to represent gender, the second component to represent how royal it is, then the age and then whether it's a food, and so on."

for the learned word embedding using one of the algorithms that we've seen like the GloVe algorithm, what happens is, you cannot guarantee that the individual components of the embeddings are interpretable.

3.1 Reason:

e.g right: there is some feature space with the first axis is gender and the second axis is royal. What you can do is guarantee that the first axis of the embedding vector is aligned with this axis of meaning, of gender, royal, age and food.

And in particular, the learning algorithm might choose right e_w^1 to be axis of the first dimension.

So given maybe a context of words, the first dimension might be right e_w^1 and the second dimension might be e_w^2 . (it might not even be orthogonal-正交), maybe it'll be a second non-orthogonal axis, could be the second component of the word embeddings you actually learn.

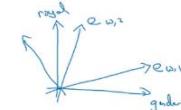
$$(\theta_i^T e_j + b_i - b'_j - \log X_{ij})^2$$

A note on the featurization view of word embeddings

	Man	Woman	King	Queen	
Gender	-1	1	-0.95	0.97	
Royal	0.01	0.02	0.93	0.95	
Age	0.03	0.02	0.70	0.69	
Food	0.09	0.01	0.02	0.01	

$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij})(\theta_i^T e_j + b_i - b'_j - \log X_{ij})^2$$

$$(\theta_i^T e_j + b_i - b'_j - \log X_{ij})^2$$



网易云

Andrew Ng

3.2 Proof: Glove minimize object: term: $(\theta_{-j})^T e_j$:

given invertible matrix A, then θ_{-i}, e_j could be easily replaced : $(A\theta_{-i})^T(A^{-1}e_j) = (\theta_{-i})^T e_j$

This is a brief proof that shows, with an algorithm like this, you can't guarantee that the axis used to represent the features will be well-aligned with what might be easily humanly interpretable axis. In particular, the first feature might be a combination of gender, and royal, and age, and food, and cost, and size, and all the other features. It's very difficult to look at individual components, individual rows of the embedding matrix and assign the human interpretation to that.

3.3 Note:

But despite this type of linear transformation $(A\theta_{-i})^T(A^{-1}e_j)$, the parallelogram map that we worked out when describing analogies, that still works. So, despite this potentially arbitrary linear transformation of the features, you end up learning the parallelogram map for figure analogies still works.

5.2-9 Sentiment classification_NLP and Word Embeddings

1. Sentiment classification problem:

1.1 Sentiment classification definition:

Sentiment classification is the task: looking at a piece of text and telling if someone likes or dislikes the thing they're talking about. It is one of the most important building blocks in NLP and is used in many applications.

1.2. Challenge:

One of the challenges of sentiment classification is might not have a huge label training set for it.

But with word embeddings, you able to build good sentiment classifiers even with only modest-size label training sets.

(input to neural network for rating is input sentence words embedding --> with trained embedding matrix, get to generalize neural network for rating, from mapping small training sentence words embedding to rating, to mapping huge, other dev/test sentence word embedding to rating, as words embedding relation between small training sentence and huge dev/test sentence have learned well by word embedding matrix E.

1.3: Example:

Input X: is a piece of text

Output Y: predict what is the sentiment, such as the star rating of, like a restaurant review.

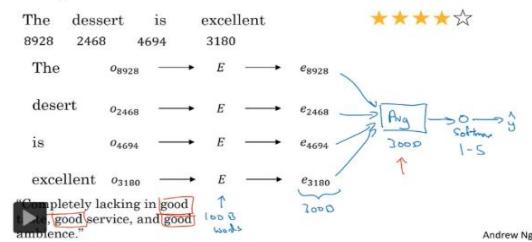
e.g: someone says, "The dessert is excellent" and they give it a four-star review, "Service was quite slow" two-star review, etc.

1.4 Summarize

If train a algorithm to map from X to Y based on a label data set like this, then could use this algorithm to monitor comments that people are saying about maybe a restaurant, it can look just a piece of text and figure out how positive or negative is the sentiment toward restaurant.

So one of the challenges of sentiment classification is you might not have a huge label data set. So for sentimental classification task, training sets with maybe anywhere from 10,000 to maybe 100,000 words would not be uncommon. Sometimes even smaller than 10,000 words. And word embeddings can help you to much better understand especially when you have a small training set.

Simple sentiment classification model



Andrew Ng

2. Simple model_sentiemnt classification

2.1 Example

Input sentence: "dessert is excellent"

output: labeled : 4 star

object: build a classifier to map it to the output Y

Look up those training sentence words in dictionary, and build a classifier to map it to the output Y that this was four stars.

2.2: Process

>1. Generate embedding vectors for all words of input sentence

>>>1. Take input sentence words and look up the one-hot vector.

>>>2. a one-hot vector multiplied by the embedding matrix E, get embedding vectors for each words

Embedding Matrix E can learn from a much larger text corpus, like a billion words or a hundred billion words, and use that to extract out the embedding vector for the word "the "dessert" is" "excellent".

If embedding matrix E was trained on a very large data set, then this allows you to take a lot of knowledge even from infrequent words and apply them to your problem, even words that weren't in your labeled training set.

>2. Build a classifier_simple one: sum/average embedding vectors and feed into softmax classifier

one simple way to build a classifier: take input sentence embedding vectors, then just sum or average them.

e.g: given a 300-dimensional feature vector , and pass it to a soft-max classifier which then outputs Y-hat: output the probabilities of the five possible outcomes from one-star up to five-star.

2.2: Process

>2. Build a classifier_simple one: sum/average embedding vectors

Note: simple mode-sum/average embedding vector

1.advantage: by using the average operation, this particular algorithm works for reviews/sentence that are short or long .

Because even if a review/sentence that is 100 words long, can just sum or average all the feature vectors for all hundred words and gives a 300-dimensional feature representation, that can then pass into sentiment classifier.

So this average will work decently well. averages or sum the meanings of all the words in your example sentence.

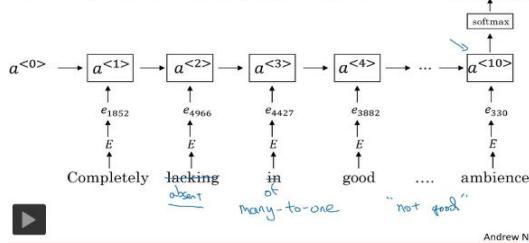
2. Disadvantage: ignores word order

So one of the problems with this algorithm is it ignores word order.

e.g: sentence "Completely lacking in good taste, good service, and good ambiance".

As the word good appears a lot, the algorithm like this will ignore word order and just sums averages all of the embeddings for the different words, then you end up having a lot of the representation of good in your final feature vector and then classifier will probably think this is a good review while it is a one-star review.

RNN for sentiment classification_many to one architecture



Andrew Ng

3. RNN for sentiment classification_many to one architecture

Input sentence: "Completely lacking in good taste, good service, and good ambiance"

3.1 Process:

>1. Generate one-hot vector for each word of input sentence O_i

>2. Generate embedding vector for each word of input sentence: $e_{-i} = E \cdot O_i$

>3. Feed all embedding vectors of input sentence words, to RNN.

RNN job is to compute the representation of input sentence at the last time step (e.g: right a<10>) that allows you to predict y^<10>

3.2 RNN model Property/advantage:

>1. This RNN is a many-to-one RNN architectur.

>2. And with an algorithm like this, it will be much better at taking word sequence into account :

e.g: realize "things are lacking in good taste" is a negative review unlike the above sum/aver algorithm, which just sums everything together into a big-word vector mush and doesn't realize that "lack good" has a very different meaning than the words "good" and so on.

So train this RNN algorithm, will end up with a pretty decent sentiment classification algorithm.

3.3 . Note: word embedding matrix E

As word embeddings can be trained from a much larger data set, this will do a better job generalizing to maybe even new words now that do not see in training set.

e.g: sentence "Completely absent of good taste, good service, and good ambiance", then even if the word "absent" is not in label training set, if it was in your 1 billion or 100 billion word corpus used to train the word embeddings, it might still get this right and generalize much better even to words that were in the training set used to train the word embeddings but not in the label training set for specifically the sentiment classification problem.

So that's it for sentiment classification, and I hope this gives you a sense of how once you've learned or downloaded from online a word embedding, this allows you to quite quickly build pretty effective NLP systems.

5.2-10: Debiasing word embedding _NLP and word embedding

Machine learning and AI algorithms are increasingly trusted to help make, extremely important decisions. And so we like to make sure that as much as possible that they're free of undesirable forms of bias, such as gender bias, ethnicity bias and so on.

This video will show some of the ideas for diminishing or eliminating these forms of bias in word embeddings.

1. The problem of bias in word embedding

1.1 Bias definition here: gender, ethnicity, sexual orientation bias.

That's a different sense of bias than is typically used in the technical discussion on machine learning.

1.2 Bias problem in word embedding:

gender stereotype problem:

E.G1: Word embeddings can learn analogies like man is to woman as king is to queen. But if ask it, man is to computer programmer as woman is to what: a horrifying result where a learned word embedding might output: Man to Computer Programmer as Woman to Homemaker.

It'd be much more preferable to have algorithm output man is to computer programmer as a woman is to computer programmer.

e.g 2: Father to Doctor as Mother is to what? And the really unfortunate result is that some learned word embeddings would output Mother:Nurse

This enforces a very unhealthy gender stereotype.

Word embeddings can reflect the gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model/word embedding matrix.

The problem of bias in word embeddings

Man:Woman as King:Queen

Man:Computer_Programmer as Woman:Homemaker X

Father:Doctor as Mother:Nurse X

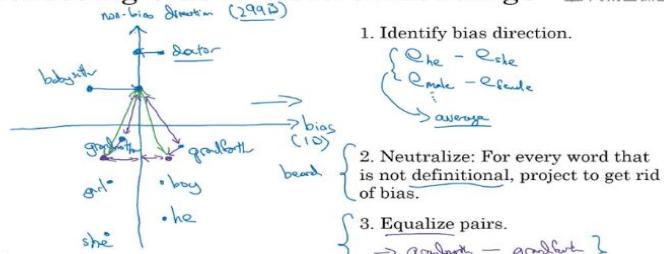
Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model.



[Boulkabasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings] ↵

Addressing bias in word embeddings

网易云课堂



Andrew Ng

2. Addressing bias in word embedding

2.1 Process

>0: Assuming already learned a word embedding: like right pic.: draw position in feature space of word: babysitter, doctor, , grandmother , grandfather, gir, boy , she, he.

>>1. first : identify bias direction:

Identify the direction corresponding to a particular bias we want to reduce or eliminate.

e.g: focus on gender bias (note: these ideas are applicable to all of the other types of bias that mentioned above).

>>>1. Identify bias direction method:

e.g: for the case of gender, take the embedding vector for definitional words and subtract several pairs of them:

e_he - e_girl

e_male - e_female

and take a few of these differences and basically average them. And this will allow you to figure out in this case the gender direction/ the bias direction- e.g right pic, horizontal direction.

Non-bias direction:

is the direction unrelated to the particular bias we're trying to address.

e.g: in this case, think of the gender bias direction as a 1D subspace whereas a non-bias direction, this will be 299-dimensional subspace.

>>>2. Note:

here has simplified the description a little bit in the original paper. The bias direction can be higher than 1-dimensional, and rather than take an average, it's actually found using a more complicated algorithm called a SVD, a singular value decomposition (奇异值分解)- principle component analysis, it uses ideas similar to the pc (主成分分析(pc)) or the principle component analysis algorithm.

>3. Equalize pairs

Equalization for pairs of words that is definitional.

e.g: such as grandmother and grandfather, or girl and boy, where you want the only difference in their embedding to be the gender.

>>>1. Motivation for equalizing definitional words:

e.g. right pic: the distance / similarity, between babysitter and grandmother is actually smaller than the distance between babysitter and grandfather. And so this maybe reinforces an unhealthy, or maybe undesirable, bias that grandmothers end up babysitting more than grandfathers.

So in equalization step, will make sure that words of definitional , like grandmother and grandfather, are both exactly the same similarity/ the same distance, from words that should be gender neutral-projected on non-bias direction already, such as babysitter or such as doctor.

>>>2. method:

There are a few linear algebra steps for that. But what it will basically do is move grandmother and grandfather to a pair of points (new positions) that are equidistant-(等轴距) from non-bias axis (like right vertical axle). And so the effect of that is that now the distance between babysitter, compared to these two words, will be exactly the same.

>>>3. Decide what word to neutralize:

In general, there are many pairs of definitional words like this grandmother-grandfather, boy-girl, sorority-fraternity, girlhood-boyhood, sister-brother, niece-nephew, daughter-son, that you might want to carry out through this equalization step.

e.g: the word doctor seems like a word you should neutralize to make it non-gender-specific or non-ethnicity-specific. Whereas the words grandmother and grandfather should not be made non-gender-specific.

The paper authors did is train a classifier to try to figure out what words are definitional: should be gender-specific and what words should not be. and it turns out that most words in the English language are not definitional-means gender is not part of the definition. and it's such a relatively small subset of words like this, grandmother-grandfather, girl-boy, sorority-fraternity, and so on that should not be neutralized.

So a linear classifier can tell what words to pass through the neutralization step to project out this bias direction, to project it on to this essentially 299-dimensional subspace. And the number of pairs want to equalize is actually also relatively small (at least for the gender example), it is quite feasible to hand-pick most of the pairs you want to equalize.

The full algorithm is a bit more complicated than present it here.

3. Summarize:

Reducing or eliminating bias of our learning algorithms is a very important problem because these algorithms are being asked to help with or to make more and more important decisions in society. In this video I shared just one set of ideas for how to go about trying to address this problem, but this is still a very much an ongoing area of active research by many researchers.

(note: after eliminatin bias: neutralize every non-definational words and equalize every definational words, the parallel relationship still be kept (linear transformation or not (non-linear transformation)))??

1.3: Bias influence on real life

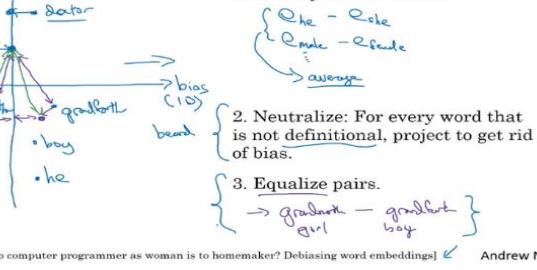
As machine learning algorithms are being used to make very important decisions. They're influencing everything ranging from college admissions, to the way people find jobs, to loan applications, whether your application for a loan gets approved, to in the criminal justice system, even sentencing guidelines. Learning algorithms are making very important decisions and so it's important that we try to change learning algorithms to diminish as much as is possible, or, ideally, eliminate these types of undesirable biases

In the case of word embeddings, it can pick up the biases of the text which used to train the model and so the biases they pick up will tend to reflect the biases in text as is written by people.

In this video will share one example of a set of ideas due to the paper referenced at the bottom by Bolukbasi and others on reducing the bias in word embeddings.

Addressing bias in word embeddings

网易云课堂



Andrew Ng

>2. second: neutralization(中和) : for every word that's not definitional (定义不明确是否是例如女性, 男性), project it to get rid of bias.

>>>1. word of definitional: there are some words that intrinsically capture gender ((定义本身就像说明了性别)): eg: grandmother, grandfather, girl, boy, she, he, a gender is intrinsic in the definition.

>>>2. word not definitional: e.g: word like doctor and babysitter, etc. that we want to be gender neutral.

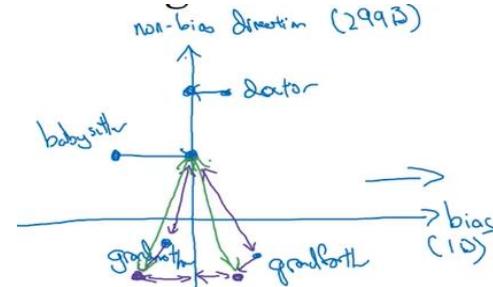
in the more general case, u might want words like doctor or babysitter to be ethnicity neutral or sexual orientation neutral, and so on, but we'll just use gender as the illustrating example here.

But so for every word that is not definitional, this basically means not words like grandmother and grandfather, which really have a very legitimate gender component, because, by definition, grandmothers are female, and grandfathers are male.

>>>3. Projection:

So for words like doctor and babysitter, just project them onto non-bias direction (right pic, vertical axis) to reduce/ eliminate their component in the bias direction- e.g. reduce their component in this horizontal direction.

怎么处理: non definitional word embedding, project on non-bias vector?



Word embedding learning algorithm:**1. Neural language model: predict the next word given previous word**> Training set:
input: previous specify window of input word embeddings
labeled output: one-hot vector of next word:

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

>Architecture: neural language model: input word embeddings feed to hidden layer 1-->softmax: (vocab size vector)
hidden layer 1: $Z_1 = W_1 \cdot [e_1, e_2, e_3, \dots] + b_1$; $a_1 = g(z_1)$
softmax: $y^a = a_2 = \text{softmax}(W_2 \cdot a_1 + b_2)$ >Parameter:
word embeddign matrix E, weight w1, w2, bias b1,b2.

$$J(\hat{y}, y) = - \sum_{i=1}^{10,000} y_i \log \hat{y}_i$$

note: architecture actually is one time step of sequence language model.

Sequence language model : many to one RNN architecture:

Training set: sentence example + labeled y (O_t-one hot vector)= [y<1>,..y<t>,..y<Ty>]

Input: a_{t-1} , $x_{t-1} = y_{t-1}$ output labeled: target word [y<1>,..y<t>,..y<Ty>], y_{t-1} is one-hot vector
model prediction y^a : softmax function for each word in vocab: probability of each word in vocab being the target word: $p(t|y_{t-1}, y_{t-2}, \dots, y_{t-1})$.Loss: for one time step: $-\sum_{i=1}^{10,000} \log(y^a_{t-1})$ --> maxmize y^a_{t-1} : i: the ith postion in vocab, labeled as 1**Word embedding learning algorithm:****2. skip gram model_V2V:- neural softmax function**

one to one RNN architecture: neural network-simplily just softmax function

>>> Training set: pair of word: context + target

input: context word embedding: e_c

output labeled: target word one-hot vector

output-softmax- vocab size category: probability of each word in vocab being the target word:

>>> Architecture: just softmax function: $y^a_t = \exp(\theta T_t \cdot e_c + b_t) / \sum \exp(\theta T_t \cdot e_c + b_t)$ over t-all words in vocab.

>>>Parameter:

word embeddings E, each vocab word t weight θT_t Loss function:-softmax loss: for one training sentence: sum $y(k) \cdot \log(y^a(k))$ over k-all vocab words.**3. Skip-gram mdoel_- negative sampling: nerual network: just 'softmax'-sigmoid function for each word in vocab**

>Training set: one example: K+1 pair: 1 positive context-target (from training sentence), k negative context-target (from vocab)

input: context word embedding e_c output labeled: vecotor like one-hot: but only labeled k+1 element: 1 for positive target word, k 0 for negative target word.
model prediction output: sigmoid function for each word in vocab-->vocab size binary classifier> Architecute: vocab size bianry classifier: for each word in vocab: $y^a_t = \text{sigmoid}(\theta T_t \cdot e_c + b_t)$.

> Cost: only compute k+1 element loss of predicted output

note: while in backprop, only use K+1 word (labeled 1/0 in output y) in each iteration, not go through all vocab word. --> in the iteration : only update k+1 words θT_t & word embedding e_c

save computation/itration.

4. Glove: neural network linear function: $\theta T \cdot e_c$ > training set: pair word i,j+ relation x_{ij} (i occure times in some window of j in all training sentece), i,j word in vocab .input: context word embedding e_c output labeled: x_{ij} model output prediction: $\theta T \cdot e_c$ >Architecture: line regression: $\theta T \cdot e_c$ > cost function: minimize $[\theta T_i \cdot e_j - \log(x_{ij})]$ sum over i (one trainig example cost, word j), sum over j (all training example)

-->One iteration: sum over vocab size twice-->more computation cost than softmax function??

(note: parameter $\theta T_i \cdot e_j$ play same role if X_{ij} is symmetric)**Word embeddign motivation: transfer learing: for new tasks have small labelle training set.**As word embeddings can be trained from a much larger data set, this will do a better job generalizing to maybe even new words now that do not see in training set.
as in the huge word corpus used to train the word embeddings, it might generalize much better even to words that: were in the training set used to train the word embeddings but not in the label training set for specifically like 'the sentiment classification' problem.

5-3 Sequence to Sequence Models

5.3-1 Basic models_Sequence to sequence models

In this week, you'll hear about sequence to sequence models, which are useful for everything from machine translation to speech recognition. Let's start with the basic models, and then later this week, you'll hear about beam search, the attention model, and we will wrap up the discussion of models for audio data like speech.

1. Sequence to sequence model:

1.1 Example:

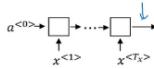
Training set:

Input: "a French sentence like Jane visite l'Afrique Septembre";
use $x^{<1>} \dots x^{<T_x>}$ as representations of input sequence

Labeled output: translate it to the English sentence, "Jane is visiting Africa in September".
use $y^{<1>} \dots y^{<T_y>}$ as representations of output sequence

1.2 Training a neural network

Note: ideas here are mainly from these two papers in right pic.



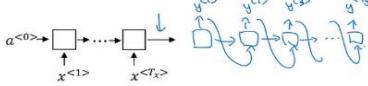
Process:

>>1. First have an encoder network - outputs a vector that represents the input sentence.

>>2. feeding the input French words one word at a time.

>>3. After ingesting the input sequence the RNN then outputs a vector C that represents the input sentence.

(when input sentence is very long, hard for this feature vector C to remember all input features, --hard to provide enough info. for decoder to translate out)



>2. Build a decoded network:

Takes the encoding output as input, and then can be trained to output the translation one word at a time. $y^{<1>} \dots y^{<T_y>}$.

Eventually, it output the end of sequence token, upon which the decoder stops.

Note: as usual, could take the generated output $y^{<t>}$ and feed them to the next step in sequence like what we were doing before when synthesizing text using sampling language model.

2. Image captioning

An architecture very similar to this-machine translation also works for image captioning.

2.1 Example:-right pic.

Input: Given an image
Output: output captions automatically like "a cat sitting on a chair"

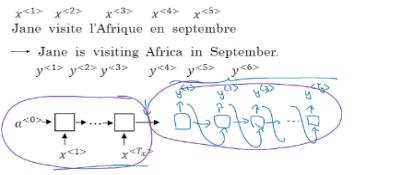
1.3 Summarize:

One of the most remarkable recent results in deep learning is that this model works: Given enough pairs of French and English sentences, if you train a model to input a French sentence and output the corresponding English translation, this will actually work decently well.

This model simply uses an encoding network whose job it is to find an encoding of the input French sentence, and then use a decoding network to then generate the corresponding English translation.

(note: input $x^{<t>}$ is one-hot vector? as machine translation have enough labeled training set, no need word embedding?)

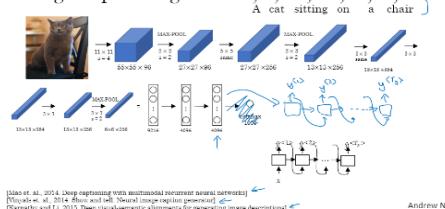
Sequence to sequence model



[Sutskever et al., 2014. Sequence to sequence learning with neural networks] ↗

[Cho et al., 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation] ↗

Image captioning



2.2 Training algorithm

Algorithm object: train network to input an image and output caption(like above phrase)

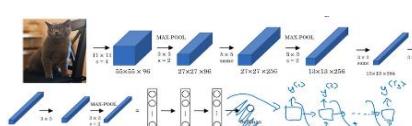
Process:

>2. Decoder network: take feature vector of input image as input, and then output caption

>>1. Feed feature vector of input image to an RNN:

>>2. RNN will generate the caption one word at a time.
Similar to what we saw with machine translation, translating from French to the English, you can now input a feature vector describing the inputs and then have it generate an output set of words, one word at a time.

This actually works pretty well for image captioning, especially if the caption you want to generate is not too long.



2.2 Training algorithm

Algorithm object: train network to input an image and output caption(like above phrase)

Process:

>1. Encode network: modified on AlexNet, output feature vector representing input image

>>1. Use AlexNet neural network:

Earlier course on the convnet have shown how you can input an image into a convolutional network, maybe a pre-trained AlexNet, and then have that learn an encoding or learn a set of features of the input image. This is actually the AlexNet architecture.

>>2. Get rid of AlexNet final softmax unit, the pre-trained AlexNet can give a multi (like 4,096) dimensional feature vector of which to represent this picture of a cat.

This pre-trained network can be the encoded network for the input image and now have a multi (like 4,096)-dimensional vector that represents the image.

(note: or the Siamese network: used for face verification??)

3. Summarize

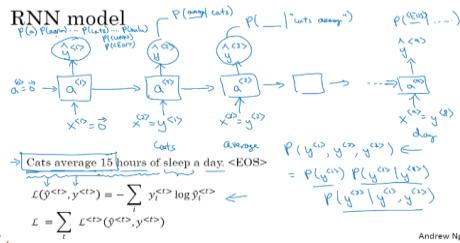
Now have seen how a basic sequence to sequence model works. How basic image to sequence- image capturing , model works. But there are some differences between how you'll run a model like this, generating a sequence compared to how you were synthesizing novel text using a language model. One of the key differences is you don't want a randomly chosen translation: may be want the most likely translation or don't want to randomly choose in caption, might want the best caption and most likely caption. Let's see in the next video how you go about generating that.

5.3-2: Picking the most likely sentence_seq 2 seq models.

There are some similarities between the sequence to sequence machine translation model and the language models that you have worked within the first week of this course, but there are some significant differences as well.

1. Machine translation as building a conditional language model

Can think of machine translation as building a conditional language model.



1.1: Language modeling

Training set: one sentence $y = y^{<1>} \dots y^{<T_y>}$

Input:

activation $a^{<0>} = 0$ vector

input in each time step: $x^{<t>} = \text{labeled } y^{<t-1>}$

(sampled language model: $x^{<t>} = \text{sampled } y^{<t-1>}$,

output in each time step-labeled: $y^{<t>} : \text{one-hot vector}$

output in each time step -predicted: soft max function:

$y^{<t>} | y^{<1>} \dots y^{<t-1>} :$ given previous labeled/ sampled words, the probability of the next word is t (in vocab).

Language modeling function:(right pic. was the network we had built in the first week);

>1. Allows to estimate the probability of a sentence $P(y^{<1>} \dots y^{<T_y>}) = p(y^{<1>}) * p(y^{<2>} | y^{<1>}) \dots p(y^{<T_y>} | y^{<1:T_y-1>})$

>2. Can also use to generate novel /new sentences:

in this example let input $x^{<t>} = \text{labeled } y^{<1>} \dots y^{<t-1>} :$ just to clean up for the slide,

1. Machine translation as building a conditional language model

Can think of machine translation as building a conditional language model.

1.2: Machine translation

Model structure: looks as right.

>1. Have encode network

Right pic. encoded network in green and the decoded network in purple.

>2. Decoded network (in purple)

The decoded network looks pretty much identical to the (sampled) language model.

Function:

Modeling the probability of, like the output English translation, conditions on some input French sentence: $P(y^{<1>} \dots y^{<T_y>} | x^{<1>} \dots x^{<T_x>})$

$$P(y^{<1>} \dots y^{<T_y>} | x^{<1>} \dots x^{<T_x>}) = \prod_{t=1}^{T_y} p(y^{<t>} | x^{<1:T_x>})$$

1.3 Compare machine translation vs language model:

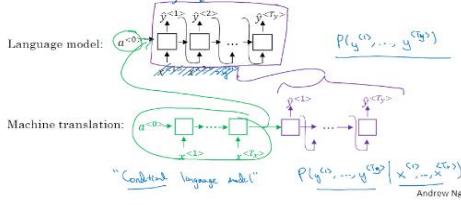
> 1. Difference: starting activation

Machine translation model is very similar to the language model, except for **starting activation**: Instead of always starting along with the vector of all zeros $a^{<0>} = \langle 0, 0, \dots, 0 \rangle$, it has an encoded network as input.

Machine translation encoded network figures out some representation for the input sentence, and it takes that input sentence and starts off the decoded network with representation of the input sentence rather than with the representation of all zeros.

-->**that is the reason machine translation a conditional language model**: instead of **modeling the probability of any sentence**, it is now **modeling the probability of, like, the output English translation, conditions on some input French sentence**. e.g.: trying to estimate the probability of an English translation: what's the chance that the translation is "Jane is visiting Africa in September," but conditions on the input French censors like, "Jane visite l'Afrique en septembre."

Machine translation as building a conditional language model



2. Find the most likely translation

Jane visite l'Afrique en septembre. $P(y^{<1>}, \dots, y^{<T_y>} | x)$

2.1: Example

Apply machine translation model to actually translate a sentence from French into English.

Input: right: French sentence: "Jane visite l'Afrique en septembre."

Model Output: the probability of different English translations for input sentence

2.2: Finding best output sentence

>1. **Attention:** Do not want to sample outputs at random.

If sample words from this distribution, $p(y^{<1>} | x)p(y^{<2>} | x) = p(y^{<1>} | x) * p(y^{<2>} | x, y^{<1>}) * p(y^{<3>} | x, y^{<1>}, y^{<2>}) \dots$
e.g.: maybe one time get a pretty good translation, "Jane is visiting Africa in September." But, maybe another time get a different translation, "Jane is going to be visiting Africa in September." Which sounds a little awkward but is not a terrible translation, just not the best one. And sometimes, just by chance, get others: "In September, Jane will visit Africa."

>2. **Object:** of while using model $p(y^{|}x)$ for machine translation: find output sentence y^{\wedge} that maximize probability $p(y^{\wedge}|x)$:
So when using this model $p(y^{\wedge}|x)$ for machine translation, do not sample at random from this distribution, but instead try to **find the English sentence y^{\wedge} that maximizes that conditional probability $p(y^{\wedge}|x)$** .

>3. **Method of find best y^{\wedge} :** come up an algorithm to finding best output sentence y^{\wedge}

So in developing a machine translation system, one of the things you need to do is **come up with an algorithm that can actually find the value of y^{\wedge} that maximizes probability: $p(y^{\wedge}<1>, \dots, y^{\wedge}<T_y> | x)$**

The most common algorithm for doing this is called **beam search**, and it's something you'll see in the next video.

3. Why not a greedy search

3.1 Definition:

Greedy search algorithm : generate the first word by just picking whatever is the most likely first word according to the conditional language model. Feed this chosen first word to machine translation model and then pick whatever is the second word that seems most likely, then pick the third word that seems most likely.

3.2 Greedy search problem

Machine translation model object: to pick the entire sequence of words, $y^{<1>}, y^{<2>}, \dots, y^{<T_y>}$, that **maximizes the joint probability of that whole thing**: $P(y^{<1>}, y^{<2>}, \dots, y^{<T_y>} | x)$

It turns out that the greedy approach, could not get maximizes the joint probability.

Greedy search just pick the best first word, and then, after having picked the best first word, try to pick the best second word, and then, after that, try to pick the best third word, doesn't really work.

Example: greedy search problem

have two output translation:

1: "Jane is **visiting** Africa in Septemeber."

2: "Jane is **going** to be visiting Africa in September"

The first one is a better translation, so expect machine translation model will have probability $p(y|x)$ is higher for the first sentence. But if we use greedy search algorithm, after picked "Jane is" as the first two words as "going" is a more common English word, probably the chance of "Jane is going," given the French input, this might actually be higher than the chance of "Jane is visiting," given the French sentence: $p(\text{going} | \text{Jane is, X}) > p(\text{visiting} | \text{Jane is, X}) \rightarrow P(\text{Jane is going}) > P(\text{Jane is visiting})$

4. Summarize

This video show how machine translation can be posed as a conditional language modeling problem. But one major difference between this and the earlier language modeling problems is rather than wanting to generate a sentence at random, you may want to **try to find the most likely English sentence**, most likely English translation. But the set of all English sentences of a certain length is too large to exhaustively enumerate. So, we have to resort to a search algorithm. So, with that, let's go onto the next video where you'll learn about beam search algorithm.

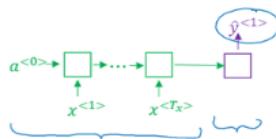
5.3-3 Beam search_Seq 2 Seq models

In the last video, you remember how for machine translation given an input French sentence, you don't want to output a random English translation, you want to output the best and the most likely English translation. The same is also true for speech recognition where given an input audio clip, you don't want to output a random text transcript of that audio, you want to **output the best, maybe the most likely, text transcript**. Beam search is the most widely used algorithm to do this.

1. Beam Search algorithm

1.1: Example_translation:

Input: French sentence, "Jane, visite l'Afrique en Septembre"
Output: expected being translated into, "Jane, visits Africa in Septemb



1.2 Process_beam search

> 0: Pre-process: generate $p(y^{<1>} | x)$

input French sentence through this encoder network and then this first step will then decode the network $y^{<1>}$, this is a softmax output overall vocab size like 10,000possibilities.

>2. Difference: 'decoder' inpute

Language model: $x^{<t>} = \text{labeled } y^{<t-1>}$: inpute is previous word

yet sampling language model inpute is the sampled previous generated word $y^{<t-1>}$.

machine tranlation: decoder inpute in the word generated in previous step: $y^{<t-1>} - \text{same as sampled language mode}$

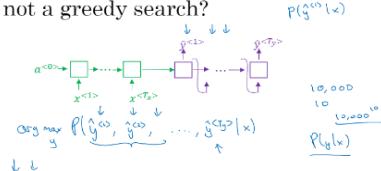
Finding the most likely translation

Jane visite l'Afrique en septembre. $P(y^{<1>}, \dots, y^{<T_y>} | x)$

- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September, Jane will visit Africa.
- Her African friend welcomed Jane in September.

$$\arg \max_{y^{<1>} \dots y^{<T_y>}} P(y^{<1>}, \dots, y^{<T_y>} | x)$$

Why not a greedy search?



- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.

$$P(\text{Jane is going} | x) > P(\text{Jane is visiting} | x)$$

Andrew

3.2: Greedy search problem

Example: greedy search problem

-->

If just pick the third word based on whatever maximizes the probability of just the first three words, might end up choosing the second translation, and ultimately ends up resulting in a less optimal sentence as measured by this model for $p(y|x)$.

Output sequence word choose

Above is an example of a broader phenomenon: when try to find the sequence of words, $y^{<1>}, y^{<2>}, \dots, y^{<T_y>}$, all the way up to the final word that together maximize the probability, it's not always optimal to just pick one word at a time.

And total number of combinations of words in the English sentence is exponentially larger: Picking words from the vocabulary size, is just a huge space of possible sentences, and it's impossible to rate them all, which is why the most common thing to do is use an approximate search out of them.
e.g have 10,000 words in a dictionary and while contemplating translations that are up to ten words long, then there are $10,000^{10}$ possible sentences that are ten words long.

Approximate search algorithm: object: it try, it won't always succeed, to pick the sentence that maximizes that conditional probability $p(y|x)$.

Even though it's not guaranteed to find the value of y that maximizes this, it usually does a good enough job.

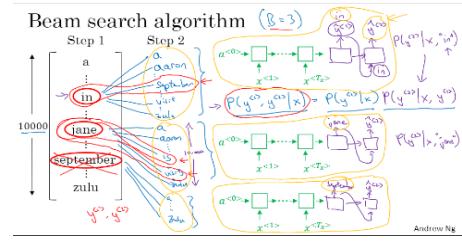
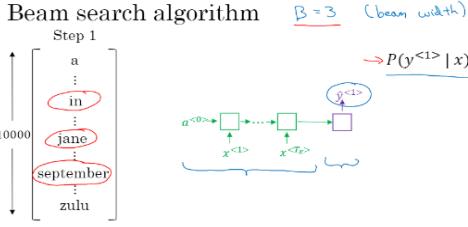
>1. First words: choose B number of $p(y^{<1>} | x)$

Take those vocab size like 10,000 possible outputs and keep in memory which were the top B for $p(y^{<1>} | x)$.

Idea: Whereas greedy search will pick only the one most likely words and move on, Beam Search instead can consider multiple B alternatives.

Parameter: B: beam width, the number of possibility. Beam search will pick at one time.

e.g: if B=3, while evaluating this probability over different choices the first words, finds that the choices 'in, Jane and September' are the most likely three possibilities for the first words in the English outputs. Then Beam search will store in computer memory that it wants to try all of three of these words, and if the beam width parameter were said differently, the beam width parameter was 10, then we keep track of not just three but of the ten, most likely possible choices for the first word.



> 2. Second word: choose B numbers of $P(y^{<1>} | x, y^{<1>} | x) = P(y^{<1>} | x) * P(y^{<2>} | x, y^{<1>} | x)$

Second step object: is to find the pair of the first and second words that is most likely $P(y^{<1>} | x, y^{<2>} | x)$, not just a second where is most likely $p(y^{<2>} | x, y^{<1>} | x)$.

>>> 1. For each of these B choices of first word:

>>>> 1. Consider what should be the second word: $p(y^{<2>} | x, y^{<1>} | x)$: have vocab size vector.

network frame used: encoder for input + decoder portion for $y^{<1>} | x, y^{<2>} | x$:

>>>> 2. Calculate likely of $P(y^{<1>} | x, y^{<2>} | x) = P(y^{<1>} | x, y^{<1>} | x) * P(y^{<1>} | x)$, get vocab size result.

for each of the B words have chosen for the first word $y^{<1>} | x$, like "in", "Jane," and "September", save away this probability $P(y^{<1>} | x)$. Then can multiply them by this second probabilities $P(y^{<2>} | x, y^{<1>} | x)$ to get the probability of the first and second words.

>>> 2. for all other B choices of first word ,do the same thing

In the second step of beam search, end up considering B x vocab size possibilities.

As: for each chosen first word, there are vocab size probabilities $P(y^{<2>} | x, y^{<1>} | x)$.

>>> 3. Pick the top B of $P(y^{<1>} | x, y^{<2>} | x)$

Evaluate all of these B x vocab size like 30000 options according to the probably the first and second words and then pick the top B. So with a cut down, these 30,000 possibilities down to the beam width-3 rounded again. that's what Beam's search would memorize away and take on to the next step beam search.

>3rd word: choose B numbers of $P(y^{<1>} | x, y^{<2>} | x, y^{<3>} | x)$

Given B pairs of first two words $y^{<1>} | x, y^{<2>} | x$, find third word, that $P(y^{<1>} | x, y^{<2>} | x, y^{<3>} | x)$ is top B.

E.G: that the most likely choices for first two words were " in September, Jane is, and Jane visits"

>>> 1. For each of these B choices of first two words:

Saved away in computer memory the probability of $P(y^{<1>} | x, y^{<2>} | x)$

>>>> 1. Consider what should be the 3rd word: $p(y^{<3>} | x, y^{<1>} | x, y^{<2>} | x)$: have vocab size vector.

network frame used: encoder for input + decoder portion for $y^{<1>} | x, y^{<2>} | x, y^{<3>} | x$:

>>>> 2. Calculate likely of $P(y^{<1>} | x, y^{<2>} | x, y^{<3>} | x) = P(y^{<3>} | x, y^{<1>} | x, y^{<2>} | x) * P(y^{<1>} | x, y^{<2>} | x)$, get vocab size result.

for each of the B pairs have chosen for the first two word $y^{<1>} | x, y^{<2>} | x$, save away this probability $P(y^{<1>} | x, y^{<2>} | x)$. Then can multiply them by this 3rd probabilities $P(y^{<3>} | x, y^{<1>} | x, y^{<2>} | x)$ to get the probability of the three words.

(note: could just pick the top B 3rd word, and multiply probability of previous chosen first two words-->get top B probabilities of three word sentence: $p(y^{<1>} | x, y^{<2>} | x, y^{<3>} | x)$)

>>> 2. for all other B choices of first two words ,do the same thing

In the 3rd step of beam search, end up considering B x vocab size possibilities.

>>> 3. Pick the top B of $P(y^{<1>} | x, y^{<2>} | x, y^{<3>} | x)$

Evaluate all of these B x vocab size like 30000 options according to the probably the first and second words and then pick the top B.

(or evaluate all B X size options-->pick top B)

2. Summary:

With a beam of B being searched considers B possibilities at a time. If the beam width =1, then this essentially becomes the greedy search algorithm.

By considering multiple possibilities say three or ten or some other number at the same time beam search will usually find a much better output sentence than greedy search. You've now seen how Beam Search works but it turns out there's some additional tips and tricks from our partners that help you to make beam search work even better. Let's go onto the next video to take a look.

5.3-4 Refinements to beam search_Seq 2 seq models.

Length normalization is a small change to the beam search algorithm that can help get much better results

1. Length normalizaiton

Beam search object: maximizing probability: $P(y^{<1>} | x, y^{<2>} | x, \dots, y^{<T>} | x) = P(y^{<1>} | x) * P(y^{<2>} | x, y^{<1>} | x) \dots * P(y^{<T>} | x, y^{<1>} | x, \dots, y^{<T-1>} | x)$

1.1. Original object formula problem:

>1.1 Numerical underflow (数值下溢出)-->log.

>>> 1. Problem:

Probabilities are all numbers often much less than 1, multiplying a lot of numbers less than 1 will result in a tiny number, which can result in numerical underflow: too small for the floating part representation in computer to store accurately.

So in practice, instead of maximizing this product, we will take logs.

>1.2: Tends to prefer very short translations

>>> 1. Problem

original objective: maximize: $P(y^{<1>} | x, y^{<2>} | x, \dots, y^{<T>} | x) = P(y^{<1>} | x) * P(y^{<2>} | x, y^{<1>} | x) \dots * P(y^{<T>} | x, y^{<1>} | x, \dots, y^{<T-1>} | x)$

if have a very long sentence, the probability of that sentence is going to be low, because multiplying as many terms here that are less than 1 to estimate the probability of that sentence.

If multiply all the numbers that are less than 1 together, you just tend to end up with a smaller probability.

And so this objective function has an undesirable effect, that maybe it **unnaturally tends to prefer very short translations**. It tends to prefer very short outputs.

As the probability of a short sentence is determined just by multiplying fewer of these numbers are less than 1. And so the product would just be not quite as small.

Note:

The same thing is true for the log (P). $P \leq 1, \log(P) \leq 0$. So the more terms p_i you have together, the more negative log (p) becomes.

>1.1.1 Numerical underflow (数值下溢出)-->log.

>>> 2. Solution: Log functon: $\log(P) = \log(p_1) + \log(p_2) + \dots + \log(p_{|Ty|})$

>>> 1. Multiplicaiton of p_i , turn in to sum of $\log(p_i)$.

>>> 2. Maximizing this sum of log probabilities should give the same results in terms of selecting the most likely sentence y.

And as log function is strictly monotonically increasing function, then maximizing $\log(P(y|x))$ should give the same result as maximizing $P(y|x)$.

>>> 3. By taking logs, you end up with a more numerically stable algorithm that is less prone to numerical rounding errors (四舍五入), or to really numerical underflow.

So in most implementations, keep track of the sum of logs of the probabilities rather than the protocol of probabilities.

Length normalization

$$\begin{aligned} & \text{arg max}_{y^{<1>} \dots y^{<T>}} \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>} | x, \dots, y^{<t-1>} | x) \\ & \text{arg max}_{y^{<1>} \dots y^{<T>}} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>} | x, \dots, y^{<t-1>} | x) \\ & \rightarrow \frac{1}{T_y} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>} | x, \dots, y^{<t-1>} | x) \end{aligned}$$

$$\begin{aligned} & \log p(y|x) \leftarrow \log p(y^{<1>} | x) + \log p(y^{<2>} | x, y^{<1>} | x) + \dots + \log p(y^{<T>} | x, y^{<1>} | x, \dots, y^{<T-1>} | x) \\ & \log p(y|x) \leftarrow \log p(y^{<1>} | x) + \log p(y^{<2>} | x, y^{<1>} | x) + \dots + \log p(y^{<T>} | x, y^{<1>} | x, \dots, y^{<T-1>} | x) \end{aligned}$$

Andrew Ng

>1.1.2: Tends to prefer very short translations

>>2. Solution Normalize length

Normalize length: normalize original object by the number of words- Ty translation, so it takes the average of the log (P_i) (<0) of each word, significantly reduces the penalty for outputting longer translations.
 a^x : $(0 < a < 1) \rightarrow a^x$ decrease with x increase -->
 $(p)^{1/Ty}$: increase with Ty increase

Heuristic method: $1/Ty^a$

instead of dividing by Ty , by the number of words in the output sentence, sometimes use a softer approach: Ty^a .

Parameter: a , is hyperparameter that can tune to try to get the best results.

$a = 0.7$.

If $a=1$, then completely normalizing by length.

If $a=0$, then not normalizing at all.

$a=0.7$: is somewhat in between full normalization, and no normalization

Note: using this way is a heuristic, no great theoretical justification for it, but people have found this works well.

3. Beam search discussion:

The Larger B is, the more possibilities is considering, and the better the sentence you probably find.
But the the more computationally expensive your algorithm is, as also keeping a lot more possibilities around.

3.1 Concerns to choose B:

> 1. If the beam width is very large: then consider a lot of possibilities, and tend to get a better result because you are consuming a lot of different options.

But it will be slower, the memory requirements will also grow, will also be compositionally slower.

> 2. If use a very small beam width: then may get a worse result as is keeping less possibilities in mind as the algorithm is running.
But get a result faster and the memory requirements will also be lower.

3. Beam search discussion:

3.2 Choose B:

> In practice, in production systems, it's not uncommon to see a beam width maybe around 10, and beam width of 100 would be considered very large for a production system, depending on the application.

> For research systems: where people want to squeeze out every last drop of performance in order to publish the paper with the best possible result. It's not uncommon to see people use beam widths of 1,000 or 3,000.

3.3 Note:

when B gets very large, there is often diminishing returns:

for many applications, there is a huge gain as increase beam width from like 1, which is very greedy search, to 3, to maybe 10. But the gains as you go from 1,000 to 3,000 in beam width might not be as big.

5.3-5 Error analysis on beam search_ Seq 2 Seq models

In the third course of this sequence of five courses, you saw how error analysis can help you focus your time on doing the most useful work for your project.

Now, beam search is an approximate search algorithm, also called a heuristic search algorithm. And so it doesn't always output the most likely sentence. It's only keeping track of B equals 3 or 10 or 100 top possibilities. So what if beam search makes a mistake? In this video, you'll learn how error analysis interacts with beam search and how you can figure out whether it is the beam search algorithm that's causing problems and worth spending time on. Or whether it might be your RNN model that is causing problems and worth spending time on.

1. Example translation output mistake

1.1: Example:

Dev set: input sentence: Jane visite l'Afrique en septembre

Dev set: labeled output (human translation) y^* : Jane visits Africa in September

Run beam search on learned RNN model (conditional language model): ends up translation y^* : Jane visited Africa last September

Model output is a much worse translation of the French sentence. It actually changes the meaning, so it's not a good translation.

1.2: Machine translation:

Machine translation model has two main components:

>1. RNN: Neural network model/ the sequence to sequence model = encoder + decoder

>2. Beam search algorithm: with beam width B.

Try to find which components contribute to translation error, RNN or beam search.

Note:

as it's always tempting to collect more training data that never hurts. So in similar way, it's always tempting to increase the beam width that never hurts or pretty much never hurts.

But just as getting more training data by itself might not get you to the level of performance you want. In the same way, increasing the beam width by itself might not get you to where you want to go.

2. Error analysis

To decide whether or not improving the search algorithm is a good use of your time, need to break the problem down and figure out what's actually a good use of your time.

2.1 Principle:

>1. RNN function (conditional language model): given input ,computes each output sequence word probability: $P(y^1|x)$, $P(y^2|x)$, ..., $P(y^T|x)$, ..., $P(y^{T-1}|x)$
(while training conditional language model, input x should = correct labeled y^{T-1} , yes: decoder actually is trained language model)

>2. Beam search function: try to find a value of y that gives that $\arg \max P(y^1|x), P(y^2|x), ..., P(y^T|x) = P(y^1|x)$, $P(y^2|x)$, ..., $P(y^{T-1}|x)$, $P(y^T|x)$

e.g. for a sentence, Jane visits Africa in September, RNN compute $P(\text{Jane}|x)$, $P(\text{visits}|x, \text{Jane})$, ..., $P(\text{September}|x, \text{Jane}, \dots)$

2. Implement: length normalization in Beach search

Normalized log likelihood objective: $\log [P(y^1|x), P(y^2|x), ..., P(y^T|x)] / Ty^a$

Process:

>1. Run beam search to certain steps: sentence up to certain length
Run beam search like 30 steps and consider output sentences up to length 30. if $B=3$, then algorithm will be keeping track of the top three possibilities for each of these possible sentence lengths, 1, 2, 3, 4 and so on, up to 30. all sentence output options = $B * \text{length}$

>2. look at all of the output sentences and score them using normalized log likelihood objective log(p) / Ty^a

>>> Take top B sentences $P(y^1|x), P(y^2|x), ..., P(y^T|x)$ in each step/sentence length $t=1, 2, \dots, 30$.
>>> compute modified objective function log [$P(y^1|x), P(y^2|x), ..., P(y^T|x)$] / Ty^a onto them -sentences that have seen through the beam search process.

>3. after all of these sentences validate this way, pick the one that achieves the highest value on this normalized log probability objective log [$P(y^1|x), P(y^2|x), ..., P(y^T|x)$] / Ty^a
And then that would be the final translation, your outputs.

(Note: process:

>1. trained encoder , decoder

>2. use decoder with Beam search-->get $B * \text{length}$ options

>3. use normalized log likelihood evaluate all options-->final output)

Beam search discussion

Beam width B?
 $\begin{array}{cccc} \text{large } B: \text{ better result, slower} \\ \text{small } B: \text{ worse result, faster} \\ \text{large } B \rightarrow 10, 100, 1000 \rightarrow 3000 \end{array}$

Unlike exact search algorithms like BFS (Breadth First Search) or DFS (Depth First Search), Beam Search runs faster but is not guaranteed to find exact maximum for $\arg \max_y P(y|x)$.

4. Summarize

There are exact search algorithms: computer science search algorithms like BFS, Breadth First Search, or DFS, Depth First Search.

The way to think about beam search is that, unlike those other algorithms, Beam search runs much faster but does not guarantee to find the exact maximum for this arg max that you would like to find. This is how beam search relates to those algorithms.

Beam search is a widely used algorithm in many production systems, or in many commercial systems. Now, in the circles in the sequence of courses of deep learning, we talked a lot about error analysis. It turns out, one of the most useful tools I've found is to be able to do error analysis on beam search. So you sometimes wonder, should I increase my beam width? Is my beam width working well enough? And there's some simple things you can compute to give you guidance on whether you need to work on improving your search algorithm. Let's talk about that in the next video.

Example

Jane visite l'Afrique en septembre.
→RNN →Beam Search BT

Human: Jane visits Africa in September. (y^*)
Algorithm: Jane visited Africa last September. (\hat{y}) ←
RNN computes $P(\hat{y}|x) \geq P(y^*|x)$
 $\begin{array}{c} \text{Jane} \xrightarrow{\text{RNN}} \text{visite} \xrightarrow{\text{RNN}} \text{l'Afrique} \xrightarrow{\text{RNN}} \text{en} \xrightarrow{\text{RNN}} \text{septembre} \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad x^{(4)} \quad x^{(5)} \end{array}$

Error analysis on beam search

Human: Jane visits Africa in September. (y^*)

Algorithm: Jane visited Africa last September. (\hat{y})

Case 1: $P(y^*|x) > P(\hat{y}|x) \leftarrow$

Beam search chose \hat{y} . But y^* attains higher $P(y^*|x)$.

Conclusion: Beam search is at fault.

Case 2: $P(y^*|x) \leq P(\hat{y}|x) \leftarrow$

y^* is a better translation than \hat{y} . But RNN predicted $P(y^*|x) < P(\hat{y}|x)$.

Conclusion: RNN model is at fault.

2.2 Process: break down problem

>1. using this model RNN to compute $P(y^*|x)$ and $P(y|x)$ using your RNN model.

$P(y^*|x) = P(y^* < 1> | x) \cdot P(y^* < 2> | x, y^* < 1>) \cdot P(y^* < 3> | x, y^* < 1>, y^* < 2>) \dots P(y^* < T> | x, y^* < 1>, y^* < 2> \dots y^* < T-1>)$

>2. Compare $P(y^*|x)$ and $P(y|x)$:

And then to see which of these two is bigger.

>>>1. $P(y^*|x)$ as output by the RNN model > $P(y|x)$ ----> Beam search fault

now sentence y^* actually attains a higher value for $P(y^*|x)$ than the y^* sentence. So come to conclusion beam search is failing to actually give the value of y that maximizes $P(y|x)$.

as one job that beam search had was to find the value of y that makes this really big. But it chose y^* , the y^* actually gets a much bigger value. So in this case, you could conclude that beam search is at fault.

should use separately?? compare $P(y^*|X) + P(y^*|X)$, and $\log P(y^*|X) / Ty^a + \log P(y|X) / Ty^a$
 RNN fault: $P(y^*|X) <= P(y^*|X)$
 Beam search fault: $\log P(y^*|X) / Ty^a > \log P(y|X) / Ty^a$

What if: $\log P(y^*|X) / Ty^a <= \log P(y|X) / Ty^a \& P(y^*|X) > P(y^*|X)$
 $\log P(y^*|X) / Ty^a <= \log P(y|X) / Ty^a \rightarrow$ beam search no fault
 $P(y^*|X) > P(y^*|X) \rightarrow$ RNN no fault ??

3. Do error analysis:

3.1 Process

>1. Go through the development set and find the mistakes that the algorithm made in the development set.

>2. Bread down, get conclusion for each mistake, contributor is RNN or beam search;

e.g. for the first mistakes, compute $P(y^*|x) = 2 \times 10^{-10}$. $P(y^*|x) = 1 \times 10^{-10}$. --> beam search is at fault. Then go through a second mistake, look at these probabilities, and think the model is at fault. And then go through more examples, sometimes the beam search is at fault, sometimes the model is at fault, and so on.

>3. Carry out error analysis to figure out what fraction of errors are due to beam search versus the RNN model.

3.2 Summarize

With an error analysis process like this, for every example in dev sets, where the algorithm gives a much worse output than the human translation, you can try to ascribe the error to either the search algorithm or to the objective function, or to the RNN model that generates the objective function that beam search is supposed to be maximizing.

And through this, you can try to figure out which of these two components is responsible for more errors. And only if you find that beam search is responsible for a lot of errors, then maybe is we're working hard to increase the beam width. Whereas in contrast, if you find that the RNN model is at fault, then you could do a deeper layer of analysis to try to figure out if you want to add regularization, or get more training data, or try a different network architecture, or something else. And so a lot of the techniques that in the third course in the sequence will be applicable there.

4. Summarize:

This particular error analysis process very useful whenever you have an approximate optimization algorithm, such as beam search that is working to optimize some sort of objective, some sort of cost function that is output by a learning algorithm, such as a sequence-to-sequence model or a sequence-to-sequence RNN that we've been discussing in these lectures. So with that, I hope that you'll be more efficient at making these types of models work well for your applications.

5.3-6 Bleu score_seq 2 seq models

One of the challenges of machine translation is that, given a French sentence, there could be multiple English translations that are equally good translations of that French sentence. So how do you evaluate a machine translation system if there are multiple equally good answers, unlike, image recognition where there's one right answer, just measure accuracy? If there are multiple great answers, to measure accuracy, the way this is done conventionally is through something called the BLEU score. So, in this optional video, I want to share with you, I want to give you a sense of how the BLEU score works.

1. Evaluation machine translation

1.1 Example

Input: French sentence: "Le chat est sur le tapis"

Labeled output: human generated translation: there are multiple, pretty good translations of this. So a different human might translate it differently.
 1: the cat is on the mat.
 2: there is a cat on the mat.

1.2 BLEU Score:

>1. BLEU score function: given a machine generated translation, allows to automatically compute a score that measures how good that machine translation.

>2. BLEU score intuition/Principle: so long as the machine generated translation is pretty close to any of the references provided by humans, then it will get a high BLEU score.

The intuition behind the BLEU score is we're going to look at the machine generated output and see if the types of words it generates appear in at least one of the human generated references. And so these human generated references would be provided as part of the dev set or as part of the test set.

Error analysis process

Human	Algorithm	$P(y^* x)$	$P(\hat{y} x)$	At fault?
Jane visits Africa in September.	Jane visited Africa last September.	2×10^{-10}	1×10^{-10}	(B)
...	...	—	—	(R)
...	...	—	—	(B)
...	...	—	—	(R)
...	...	—	—	1

Figures out what fraction of errors are "due to" beam search vs. RNN model

Andrew Ng

Evaluating machine translation

French: Le chat est sur le tapis.

→ Reference 1: The cat is on the mat.
 → Reference 2: There is a cat on the mat.
 → MT output: the the the the the the the

Precision: $\frac{7}{7}$ Modified precision: $\frac{2}{7}$ Count ("the")

Papineni et al., 2002, Bleu: A method for automatic evaluation of machine translation

网易云课堂

Bleu
bilingual evaluation unlabelled

Evaluating machine translation

French: Le chat est sur le tapis.

→ Reference 1: The cat is on the mat.
 → Reference 2: There is a cat on the mat.
 → MT output: the the the the the the the

Precision: $\frac{7}{7}$ Modified precision: $\frac{2}{7}$ Count ("the")

Papineni et al., 2002, Bleu: A method for automatic evaluation of machine translation

网易云课堂

Bleu
bilingual evaluation unlabelled

2. Method to measure how good machine translation output is

Example: MT output: "the the the the the the."

2.1: Basic Precision for isolated word: of the machine translation output:

>1. MT output basic precision definition: Look at each the words in the output and see if it appears in the references; measure what fraction of the words in the MT output also appear in the references.

n = MT output sentence word number;

n_i = i-th word of output sentence, appear in reference 1 or not 0;

Precision = sum n_i over all words / n

2.2. Basic precision problem:

This is not a particularly useful measure, because it seems to imply that this MT output has very high precision.

e.g. in above example, there are seven words in the machine translation output. And every one of these 7 words appears in either Reference 1 or Reference 2. So the word 'the' appears in both references. So each of these words looks like a pretty good word to include. --> So this MT output have a precision of 7 / 7.

It looks like it was a great precision.

So this is how the basic precision measure of what fraction of the words in the MT output also appear in the references. This is not a particularly useful measure, because it seems to imply that this MT output has very high precision.

2.2: Modified Precision for isolated word: of the machine translation output

>1. Modified precision definition: Look at each the words in the output and see if it appears in the references; in measure give each word credit only up to the maximum number of times it appears in each reference sentences.

n = MT output sentence word number;

n'_i = i-th word of output sentence, appear times in reference (credit only to the maximum number of times it appears in each reference sentence);

n''_i = i-th word of output sentence, appear times in per se

Precision = sum n'_i over all words / sum n''_i over all output words.

e.g:

word 'the' occurs in output 7 times, n''_i=7 (i=1);

word 'the' in Reference 1 appears twice, in Reference 2 appears just once. n'_i = 2

So, with a modified precision, it gets a score : 2/7

Modified Precision the denominator : is the sum of the count of the number of times each output word, appears per se.

Modified Precision numerator: is the sum of the max count of the number of times each output word appears in each reference.

2. Problem:

Have been looking at words in isolation.

3. Bleu score on bigrams:

In the BLEU score, you don't want to just look at isolated words. You maybe want to look at pairs of words as well.

>3.1 Definition:

Bigrams: means pairs of words appearing next to each other.

Blue score on bigrams: a portion of the final BLEU score

final BLEU score: take unigrams-single words, as well as bigrams-pairs of words into account as well as maybe even longer sequences of words, such as trigrams,-three words pairing together.

3.2 Example: MT output2: "The cat the cat on the mat". Still not a great translation, but maybe better than the last one.

3.3 Blue score computing :

>1. List all the bigrams in MT output sentence

e.g.: 'the cat', cat the, the cat - as already had that, so skip that, cat on, on the, the mat.

So these are the bigrams in the machine translation output.

>2. Count up, How many times each of these bigrams appear in MT output per se

e.g.: "The cat" appears twice, cat the appears 1, and the others all appear just once.

>3. Define the clipped count, count, and then subscript clip in reference output

Take all the bigrams, and give our algorithm credit only up to the maximum number of times that bigram appears in either Reference 1 or Reference 2.

e.g.: the cat appears a maximum of once in either of the references. So clip that count to 1. Cat the, doesn't appear in Reference 1 or Reference 2, so clip that to 0. Cat on, appears once, give it credit for once. On the appears once, give that credit for once, and the mat appears once. So these are the clipped counts.

We're taking all the counts and clipping them, really reducing them to be no more than the number of times that bigram appears in at least one of the references.

>4. modified bigram precision computation:

modified bigram precison = sum of the count clipped / sum of the count per se(MT output)

e.g.: sum of the count clipped is $1+0+1+1=4$ / sum of the count per se(MT output)= $2+1+1+1=6=4/6$

Bleu score on bigrams

Example: Reference 1: The cat is on the mat. ←

Reference 2: There is a cat on the mat. ←

MT output: The cat the cat on the mat. ←

Count	Count _{clip}	
the cat	2 ←	1 ←
cat the	1 ←	0
cat on	1 ←	1 ←
on the	1 ←	1 ←
mat	1 ←	1 ←

Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Bleu score on unigrams

网易云课堂

Example: Reference 1: The cat is on the mat.

Reference 2: There is a cat on the mat.

MT output: The cat the cat on the mat. ↑

$$P_1 = \frac{\sum_{\text{unigrams}} \text{Count}_{\text{clip}}(\text{unigram})}{\sum_{\text{unigrams}} \text{Count}(\text{unigram})}$$

$$P_n = \frac{\sum_{n\text{-grams}} \text{Count}_{\text{clip}}(n\text{-gram})}{\sum_{n\text{-grams}} \text{Count}(n\text{-gram})}$$

Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng

4. Bleu score on unigrams

4.1 Modified precision on unigrams: P1

P stands for precision and the subscript 1 here means that we're referring to unigrams, meaning we're looking at single words in isolation.

p1 denominator (分子) : sum of unigrams in the machine translation output of count, number of counts of that unigram,
p1 numerator (分子) : sum of unigrams in the machine translation output of count, number of count clip of that unigram,

4.2 Modified precision on n-gram: Pn, instead of unigram, for n-gram.

p1 denominator (分子) : sum of n-grams in the machine translation output of count, number of counts of that unigram,
p1 numerator (分子) : sum of n-grams in the machine translation output of count, number of count clip of that unigram,

And so these modified precision scores, measured on unigrams or on bigrams, or on trigrams, or even higher values of n for other n-grams. This allows you to measure the degree to which the machine translation output is similar or maybe overlaps with the references.

5. Bleu details:

Finally, let's put this together to form the final BLEU score.

Pn: the BLEU score computed on n-grams only. Also the modified precision computed on n-grams only.

5.1 Final BLEU Score: combined Bleu score: BP * exp (average (P1, P2, ...Pn))

Combining the computed P1, P2, ..., Pn together by above formula.

>1. average (P1, P2, ...Pn):

e.g.: by convention to compute one number, you compute P1, P2, P3 and P4, take the average of them

>2. exp (average):

By convention the BLEU score is defined as exponentiations (求幂) of average score, exponentiation is strictly monotonically increasing operation.

5. Bleu details:

>3. BP:BP penalty:

Added to adjust convention BLEU score.

>>>1. Motivaiton: it turns out that if MT output very short translations, it's easier to get high precision Pi. Because probably most of the words you output appear in the references.

But we don't want translations that are very short.

>>>2. Funtion: BP is an adjustment factor that penalizes translation systems that output translations that are too short.

>>>3. Formula:

If length of MT output > reference length, BP = 1, (>min(references??))

Otherwise: BP = $e^{-(\text{length of MT} / \text{length of reference})}$: $1 < \text{BP} < 2.7$, should ≤ 1 , something wrong?

Note: the details can find in this paper-right.

Bleu details

p_n = Bleu score on n-grams only

Combined Bleu score: $\text{BP} \cdot \exp\left(\frac{1}{n} \sum_{n=1}^N p_n\right)$

P_1, P_2, P_3, P_4

BP = bigram penalty

$$\text{BP} = \begin{cases} 1 & \text{if } \text{MT_output_length} > \text{reference_output_length} \\ \exp(1 - \text{MT_output_length}/\text{reference_output_length}) & \text{otherwise} \end{cases}$$

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng

Bleu details

p_n = Bleu score on n-grams only

Combined Bleu score: $\text{BP} \cdot \exp\left(\frac{1}{n} \sum_{n=1}^N p_n\right)$

P_1, P_2, P_3, P_4

BP = bigram penalty

$$\text{BP} = \begin{cases} 1 & \text{if } \text{MT_output_length} > \text{reference_output_length} \\ \exp(1 - \text{MT_output_length}/\text{reference_output_length}) & \text{otherwise} \end{cases}$$

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng

6. Summarize

1. Earlier in this set of courses, have shown the importance of having a single real number evaluation metric: it allows you to try out two ideas, see which one achieves a higher score, and then try to stick with the one that achieved the higher score.

2. So the reason the BLEU score was revolutionary for machine translation was: because this gave a pretty good, by no means perfect, but pretty good single real number evaluation metric. And so that accelerated the progress of the entire field of machine translation.

3. In practice, few people would implement a BLEU score from scratch. There are open source implementations that you can download and just use to evaluate your own system. (seems very simple algorithm, need to C/O?)

4. Today, BLEU score is used to evaluate many systems that generate text.

Such as machine translation systems, as well as the example I showed briefly earlier of image captioning systems where you would have a system, have a neural network generated image caption. And then use the BLEU score to see how much that overlaps with maybe a reference caption or multiple reference captions that were generated by people.

So the BLEU score is a useful single real number evaluation metric to use whenever you want your algorithm to generate a piece of text. And you want to see whether it has similar meaning as a reference piece of text generated by humans. This is not used for speech recognition, because in speech recognition, there's usually one ground truth. And you just use other measures to see if you got the speech transcription on pretty much, exactly word for word correct. But for things like image captioning, and multiple captions for a picture, it could be about equally good, or for machine translation there are multiple translations, but equally good. The BLEU score gives you a way to evaluate that automatically and therefore speed up your development.

(Note: BLEU: bilingual evaluation understudy: method:

1. basic precision: for isoceon word in MT: credit 1/0 if occur/not occur in ref.-->sum credit of all word in MT output sentence/ MT output sentence length

2. modified precision: for unigram/bigram/trigram word in MT: count clip/count (in MT per se)

-->final BLEU used: BP EXP (average of sum p1, ..pn) n: nigram of modified precision.)

5.3-7 Attention model intuition_seq 2 seq models

For most of this week, you've been using a Encoder-Decoder architecture for machine translation. Where one RNN reads in a sentence and then a different one RNN outputs a sentence. Attention Model make a modification to this, and it makes all this work much better. The attention idea has been one of the most influential ideas in deep learning.

1. The problem of long sequence

Example: input long sentence as right.

1.1 MT output principle: encoder + decoder:

Encoder read in the whole sentence and then memorize the whole sentences and store it in the activations conveyed here. Then the decoder network will generate the English translation.

1.2 Human translator principle:

human translator kind of , work part by part through the sentence. Because it's just really difficult to memorize the whole long sentence like that.

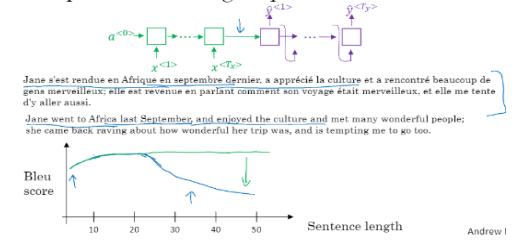
Not to first read the whole French sentence and then memorize the whole thing and then regurgitate an English sentence from scratch. Instead :

>1. Read the first part of input sentence, maybe generate part of the translation.

>2. Look at the second part, generate a few more words.

>3. Look at a few more words, generate a few more words and so on.

The problem of long sequences



1.3 MT output problem:

The Encoder-Decoder architecture above, works quite well for short sentences, so we might achieve a relatively high Bleu score, but for very long sentences, maybe longer than 30 or 40 words, the performance comes down.

1.3.1 MT output Bleu score vs sentence length-right blue curve:

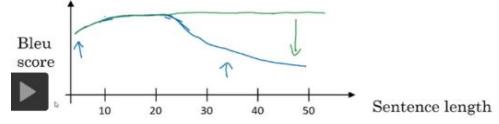
>1. Sentence is very short: Bleu score small

short sentences are just hard to translate, hard to get all the words

>2. Sentence reasonable long, Bleu score is big

>3. Sentence very long, Bleu score is small

Very long sentences, it doesn't do well because it's just difficult to get in your network to memorize a super long sentence. (encoder result C , representation of input sentence, could not remember all input sentence info. when it is too long.)



1.3.2 With Modified MT output, attention Model:

which translates maybe a bit more like humans might, looking at part of the sentence at a time, machine translation systems performance can look like this-right green curve: Bleu score is big even input sentence is very long:

as by working one part of the sentence at a time, you don't see this huge dip (gap between green curve and blue curve while sentence is very long), which is really measuring the ability of a neural network to memorize a long sentence, which maybe isn't what we most badly need a neural network to do.

(note: as long as encoder result C could remember enough info. of input sentence, than output sentence BELUE is high-->green curve,

but may not need this remembering long sentence ability to achieve this high BLEU, could modify architecture to reach this high BLEU also)

2. Attention model intuition

The Attention Model was due to Dimitri Bahdanau, Camrunc Cho, Yoshe Bengio and even though it was obviously developed for machine translation, it spread to many other application areas as well.

This is really a very influential, seminal paper in the deep learning literature.

Example:

illustrate this with a short sentence: Jane visite l'Afrique en Septembre
even though these ideas were maybe developed more for long sentences, but it'll be easier to illustrate these ideas with a simpler example.

Architecture:

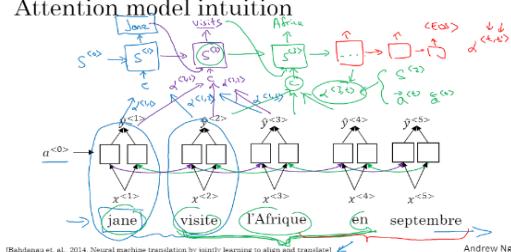
2.1. 'Encoder'-BRNN:

Generate features of each word/maybe surrounding words in input sentence word position

In this example use a modified BRNN, in order to compute some set of features for each of the input words.

BRNN: Standard BRNN outputs $y<1>, y<2>, \dots, y<Ty>$.

Modified BRNN here: standard BRNN modified here as we're not doing a word for word translation: get rid of the y on top, instead using a BRNN for each of the word/time sequence in input sentence. compute a very rich set of features about the words in the sentence and maybe surrounding words in every position.



2.2. 'Decoder'- RNN: Generate translation output words

Use another RNN to generate the English translations.

Activation Denote:

in order to avoid confusion with the activations here, use a different notation 'S' to denote the hidden state in this RNN.

Horizontal forward: \$<0>----->S<1>----->S<2>----->...----->S<Ty>,

Vertical forward: y<1>, y<2>, ...y<Ty>

Process:

>1. First word y<1>

when trying to generate this first output word y<1>, what part of the input French sentence should you be looking at.

e.g. In this example seems like should be looking primarily at this first input word 'jane', maybe a few other words close by, but don't need to be looking way at the end of the sentence.

>>> 1. Attention Model function: compute a set of attention weights and we're going to use: a<1, t>, t=1...Tx

a<1,1>; denot when generating the first words, how much should you be paying attention to this first piece of information here: feature of the first word of input sentence (which is y<1> in standard BRNN)

a<1, 2>; then we'll also come up with a second weight, which tells us what trying to compute the first work of Jane, how much attention we're paying to this second feature vector from the second word in inputs and so on and a<1,3> and so on, to a<1, Ty

>2. Denote C of all these activation together: tell what is exactly the context from denoter C that we should be paying attention to.

>3. Input C<1> to RNN unit S<1>, then try to generate the first word y<1>

That's one step of the RNN

3. Note:

1. context C<t>: the goal of the context C is: for the t-th output word, is really should capture, maybe we should be looking around surrounding (like 2) words of of t-th input word.

2. a<3,t>:

when trying to generate the third word, the amounts that this RNN step should be paying attention to the input (French) word that time t, that depends on:

>1. a<1>; The activations of the bidirectional RNN at time t, which depends on the forward prop activations and backward activations at time t

>2. s<2> and it will depend on the state from the previous steps, S<2>

These things together will influence, how much you pay attention to a specific word in the input French sentence.

4. The key intuition : 'Decoder' RNN marches forward generating one word at a time, until eventually it generates maybe the EOS and at every step, there are these attention weights a<t,t>, that tells it, when you're trying to generate the t-th output (English) word, how much should you be paying attention to the t-th input (French) words.

And this allows it on every time step to look only maybe within a local window of the input (French) sentence to pay attention to, when generating a specific output (English) word.

(note: how works actually in live:

1. need to input long sentence, then make prediction? --> could not make prediction before injecting whole sentence?

2. need to compute weight on all input sentence words' feature, input to decode-context c<t, t> also need to remember enough /compute all input words weight/info.- do not save from memory long sentence?)

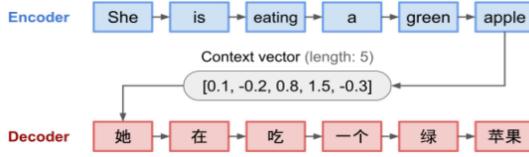
原始的 Seq2Seq 模型有什么问题

Seq2Seq 是为了解决 language modeling 问题而被提出的 (Sutskever, et al. 2014). 通俗的说, 它的目标是将输入的一个句子(source)转化成一个新的目标句子(target), 同时输入和输出的句子长度可以是任意的. Seq2Seq 的应用领域比较广泛, 有机器翻译(文字或者语音)、自动问答和自动摘要等等。。。

传统的 Seq2Seq 都拥有一个 encoder-decoder 的架构:

- Encoder 负责处理输入, 并将输入的信息进行 encode, 转化成一个固定长度的 context 向量, 这个 context 向量包含了整个输入的信息。
- Decoder 负责接收 context 向量并将其转化成输出, 早期的 Seq2Seq 模型只是利用了 encoder 最后一个状态作为 decoder 的初始状态。

Encoder 和 Decoder 通常使用 LSTM 或者 GRU。



The encoder-decoder model, translating the sentence "she is eating a green apple" to Chinese. The visualization of both encoder and decoder is unrolled in time.

将 context 设定为一个固定长度的向量有一个明显的缺点, 就是当句子长度增加之后, 它并不能很好的记住所有信息, 一个典型的栗子就是, 通常当他处理完所有的输入之后, 输入的前一部分很多信息就已经遗忘啦。。。很显然, 想用一个 context 向量来承担编码输入句子的重担, 这任务有点重。。

所有, 为了解决这个问题, Attention Mechanism (Bahdanau et al., 2015) 应运而生。

5.2-8 Attention model_seq 2 seq model

Attention model allows a neural network to pay attention to only part of an input sentence while it's generating a translation, much like a human translator might. Let's now formalize that intuition into the exact details of how you would implement an attention model.

1. Attention model

Example: input sentence-as right.

1.1: Encoder:-BRNN--->compute features on every word

>1. Encoder neural network:

Here use BRNN, or bidirectional GRU, or bidirectional LSTM to compute features on every word instead of standard BRNN output-translation.

In practice, GRUs and LSTMs are often used for this, with maybe LSTMs be more common.

>2. Denotion

$\alpha^{<1>} = (\alpha^{<1>}_f, \alpha^{<1>}_b)$: feature vector'-th word in input sentence

$\alpha^{<1>}_f$: activation in forward prop of t'-th word in input sentence

$\alpha^{<1>}_b$: activation in backward prop of t'-th word in input sentence

Forward: $\alpha^{<0>}----->\alpha^{<1>}_f----->\alpha^{<5>}_f$

backward: $\alpha^{<6>} ----->\alpha^{<5>}_b----->\alpha^{<1>}_b$

note: $\alpha^{<0>}, \alpha^{<6>}$ could set to 0 vector.

2.2. 'Decoder'- RNN: Generate translation output words

Process:

>2. For the second step of this RNN

>1. have a new hidden state S<2> and have a new set of the attention weights: a<2,1>, a<2,2>...a<2,Tx>

>2. first word generated y^<1> also feed as input

>3. C<2>: of together all a<2,1>, ...a<2, Tx>, some context that we're paying attention to a

>4, feed S<1>, y<1>, C<2> to S<2>, generated the second output word y<2>

> Third step,

>1. have a new hidden state S<3> and new context C<3> that depends on the various a<3,1>, ...a<3, Tx>.

C<3>: tells how much should we be paying attention to the different words from the input French sentence and so on.

为 Translation 而生

Attention Mechanism 起初为了在机器翻译任务中, 帮助更好的记忆长输入的情形, 和原始的只是通过 encoder 最后一个 state 传递出 context 向量不同的是, Attention 在 context 向量和 input 之间构造了多个 shortcuts (可以理解成 Attention 矩阵). 这些 shortcut 的权重相对于不同的 output, 都是不同的。

现在, 在有了 Attention Mechanism 之后, 相当于我们有了新的 context 向量, 此时, 输入和输出之间的对齐关系 (alignment relation) 可以被 context 向量学习到。

本质上来说, context 向量由以下三个部分组成。

- encoder hidden states;
- decoder hidden states;
- alignment between source and target.



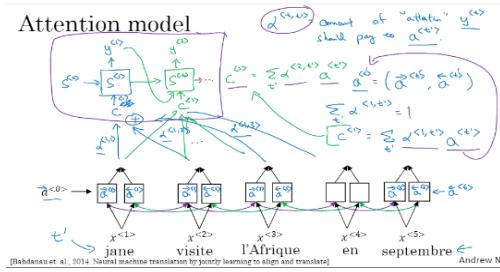
引入Attention

上面的模型使用的还是传统的 RNN, 所以在处理长文本的时候比较困难, 效果也不太理想, 这个时候就引入了注意力机制, 可以让模型只关注输入相关的部分。

加入注意力机制后, 和以前的Seq2seq有两点不同。

第一点: 以前的Seq2seq是只传递最后一个的隐藏层, 而加入Attention之后会传递所有隐藏层。

encoder

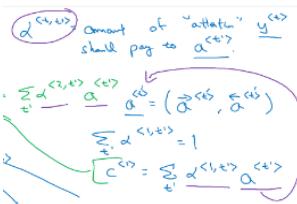


[Bahdanau et al., 2014] Neural machine translation by jointly learning to align and translate

1.2 Decoder:

$y_{<t>} = g_1(W_y S_{<t>} + b_y)$
 $S_{<t>} = g_2(W_s [y_{<t-1>} c_{<t>}] + b_s)$
 $c_{<t>} = \text{sum } a_{<t,t>} * \alpha_{<t>} \text{ over } t' \text{ (input time step)}$

$\alpha_{<t,t>} : [\alpha_{<t>}, f, \alpha_{<t>}, b]$: generated by encoder using modified BRNM
 $\alpha_{<t,t>} = \exp(e_{<t,t>}) / \text{sum}(\exp(e_{<t,t>}) \text{ over input words } t=1 \dots T_x)$
 $e_{<t,t>} = g_3(W_e [a_{<t>} s_{<t-1>}] + b_e)$



1.2.1 Neural network:

As in decoder only have our forward prop, so it's a single direction RNN with state $S_{<t>}$ to generate the translation.

1.2.2 Process:

>1. First word: $S_{<0>} + C_{<1>} \rightarrow y_{<1>}$

>>>1. Input: $S_{<0>}$, and some context $C_{<1>}$

>>>2. Context $C_{<1>} := \text{sum } a_{<1,t>} * \alpha_{<t>} \text{ over } t' \text{ all input words.}$
 depend on the attention parameters $a_{<1,1>} , a_{<1,2>} , \dots$, and so on.

>>>3. Attention parameter: $a_{<1,t>}$:

tells how much the context $C_{<1>}$ would depend on the features $\alpha_{<t>}$ gotten from the different time steps of input sentence.
 So the way we define the context C is actually be a way to sum the features from the different time steps weighted by these attention weights.

Attention weight $a_{<t,t'>}$ satisfy:

>>>>>1. All be non-negative, ≥ 0

>>>>2 sum over all time steps = 1: sum $a_{<t,t'>}$ over all t' words in input sentence = 1

So $a_{<t,t'>}$: is the amount of attention that $y_{<t>}$ should pay to a $a_{<t'>}$; in other words, when generating the t -th of the output words, how much you should be paying attention to the t' -th of input word.

So that's one step of generating the output and then at the next time step.

2. Computing attention $a_{<t,t'>}$

2.1 Definition: $a_{<t,t'>}$:

$a_{<t,t'>} : \text{the amount of attention should paid to } \alpha_{<t'>} \text{ when trying to generate the } t\text{-th words in the output translation.}$

2.2 Formula: $a_{<t,t'>} = \exp(e_{<t,t'>}) / \text{sum}(\exp(e_{<t,t'>}) \text{ over input words- } t'=1 \dots T_x)$

Note: use essentially a softmax to make sure that these weights $a_{<t,t'>}$ sum to one if sum over t' .
 So for every fix value of t , these weights sum to one if summing over t' .

2.3 Factor $e_{<t,t'>}$ computation: use a small neural network $\rightarrow a_{<t,t'>} = \frac{\exp(e_{<t,t'>})}{\sum_{t'=1}^{T_x} \exp(e_{<t,t'>})}$
 One way to compute factor $e_{<t,t'>}$ is to use a small neural networ

2.3.1 Input:

>1. $S_{<t-1>}$: neural network state from the previous time step

While trying to generate $y_{<t>}$, the hidden state from the previous step $S_{<t-1>}$ just fell into $S_{<t>}$.

$S_{<t>}$ is one hidden layer in neural network because you need to compute these $S_{<t>}$ a lot.

>2. $a_{<t>}$: feature vector of t -th input word.

2.3 Factor $e_{<t,t'>}$ computation: use a small neural network

2.3.2 Neural network intuition of $e_{<t,t'>}$:

If you want to decide how much attention to pay to the activation $\alpha_{<t'>}$, seems like it should depend the most on: what is your own hidden state activation from the previous time step.

You don't have the current state activation $s_{<t>}$ yet because of context $C_{<t>}$ feeds into $s_{<t>}$, so you haven't computed $S_{<t>}$.

But look at whatever the hidden stages of this RNN generated in the upper translation $S_{<t-1>}$ and then look at the features of each input words $x_{<t>}$. So it seems pretty natural that weight parameter $a_{<t,t'>}$ and $e_{<t,t'>}$ should depend on these two quantities $\alpha_{<t'>}, S_{<t-1>}$.

But we don't know what the function is to generate $e_{<t,t'>}$ based on $\alpha_{<t'>}, S_{<t-1>}$.
 So one thing could do is just train a very small neural network to learn whatever this function should be. And trust backpropagation, trust gradient descent to learn the right function.

It turns out that if you implemented this whole model and train it with gradient descent, the whole thing actually works.
 This little neural network does a pretty decent job telling you how much attention $y_{<t>}$ should pay to $\alpha_{<t'>}$.

Note:

what is this neural architecture: layer number, hidden units number, etc.

what's cost function of $a_{<t,t'>}$?? no need, combine with whole RNN optimizing?

3. Attention examples

3.1 Date normalization problem:

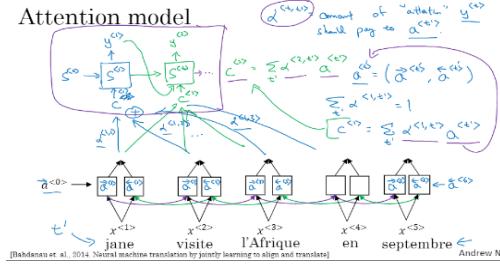
Machine translation is a very complicated problem, in the prior exercise get to implement and play of the attention while for the date normalization problem: inputting a date like "July 20th 1969" and normalizing it into standard formats like "1969-07-20".
 you can train a neural network to input dates in any of these formats and have it use an attention model to generate a normalized format for these dates.

3.2 visualize attention weight $a_{<t,t'>}$

right pic. example: a machine translation example and plotted in different colors, the magnitude of the different attention weights. could find that for the corresponding input and output words, the attention weights will tend to be high. Thus, suggesting that when it's generating a specific word in output, usually paying attention to the correct words in the input.
 All this including learning where to pay attention, when, was all learned using backpropagation with an attention model.

So that's it for the attention model really one of the most powerful ideas in deep learning

Attention model



2. Second word: $y_{<1>} + C_{<2>} (\cdots S_{<1>} \cdots) \rightarrow y_{<2>}$

Performe similarly as step 1, only now have a new set of attention weights $a_{<2,t'>}$ and they find a new way to sum, generates a new context $C_{<2>}$.

Note: now $y_{<1>}$ is also input as well as $S_{<1>}, C_{<2>}$ to generate the second word.

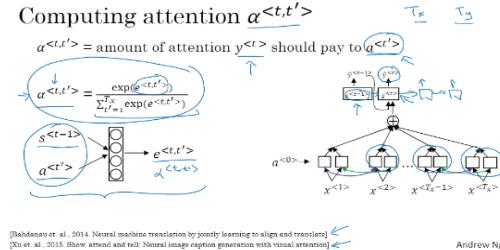
Summarize:

Using these context vectors C_1, C_2 , and so on as the input of decoder.
 The encoder network looks like a pretty standard RNN sequence with the context vectors $C_{<t>}$ as output and we can just generate the translation $y_{<t>}$ one word at a time.

Computing attention $a_{<t,t'>}$

$$a_{<t,t'>} = \text{amount of attention } y_{<t>} \text{ should pay to } \alpha_{<t'>}$$

$$\alpha_{<t,t'>} = \frac{\exp(e_{<t,t'>})}{\sum_{t'=1}^{T_x} \exp(e_{<t,t'>})}$$



2.4 Summarize:

1. As chug along generating one word at a time, this neural network actually pays attention to the right parts of the input sentence that learns all this automatically using gradient descent.

2. Disadvantage: of this algorithm: it take quadratic time / quadratic cost (三次方) to run this algorithm:

e.g: input sentence and output sentence lenght $T_x T_y$, then the total number of these attention parameters is $T_x T_y$. And so this algorithm runs in quadratic cost.

In machine translation applications neither input nor output sentences is usually that long, maybe quadratic cost is actually acceptable, and there is some research work on trying to reduce costs as well.

3. This idea has been applied to many other problems as well, like image capturing.

e.g: in the image capturing problem the task is to look at the picture and write a caption for that picture. So Look at the picture and pay attention only to parts of the picture at a time while you're writing a caption for a picture.

In paper set to the bottom by Kevin Chu, Jimmy Barr, Ryan Kiros, Kelvin Shaw, Aaron Korver, Russell Zarkutov, Virta Zemo, and Andrew Benjo they also showed that you could have a very similar architecture.

Attention examples

July 20th 1969 → 1969 - 07 - 20

23 April, 1564 → 1564 - 04 - 23



5.3-9 Speech recognition_Audio data

One of the most exciting developments were sequence-to-sequence models has been the rise of very accurate speech recognition. here to give you a sense of how these sequence-to-sequence models are applied to audio data, such as the speech.

1. Speech recognition problem

1.1 Definition:

Speech recognition problem: given an audio clip x , to automatically find a text transcript y .

1.2 Training set

>1. Input: audio clip:

audio clip: saying "the quick brown fox"

plot it looks right pic (grey pic): the horizontal axis is time, vertical is air pressure;

Note: sound principle

Microphone measures minuscule changes in air pressure, and the way of hearing my voice is: your ear detects little changes in air pressure, generated either by speakers or a headset.

And some audio clips like this plots with the air pressure against time.

>2. Output:

transcript. Hopefully, a speech recognition algorithm can input that audio clip and output that transcript "the quick brown fox".

2. Speech recognition model history:

2.1. Build phonemes (音节) cell:

speech recognition systems used to be built using phonemes- hand-engineered basic units of cells:

e.g. "the quick brown fox" represented as phonemes: simplify for illustration: Quick, has a "ku" and "wu", "ik", "k" sound.

Linguist used to write off these basic units of sound, and try to break language down to these basic units of sound. and linguists used to hypothesize that writing down audio in terms of these basic units of sound called phonemes would be the best way to do speech recognition.

2.2. End-to-end deep learning <-----large data sets

with end-to-end deep learning, phonemes representations are no longer necessary.

Instead, can build systems that input an audio clip and directly output a transcript without needing to use hand-engineered representations.
this is contribute to the much larger data sets.

Data size:

Academic data sets on speech recognition might be as a 300/ 3000 hours which is considered reasonable size.

for commercial systems, are now trains on over 10,000 hours and sometimes over a 100,000 hours of audio.

And it's really still moving to a much larger audio data sets-transcribe audio data sets both x and y , together with deep learning algorithm, has driven a lot of progress in speech recognition.

3. Speech recognition model

3.1 Attention model

encoder: like BRNN

Input: on the horizontal axis, take in different time frames of the audio input,

decoder: attention model

Have an attention model try to output the transcript like, "the quick brown fox"

3. Speech recognition model

3.2: CTC cost model for speech recognition

CTC: Connectionist Temporal Classification

due to Alex Graves, Santiago Fernandes, Faustino Gomez, and Jürgen Schmidhuber.

>1. Idea: use a neural network structured with an equal number of input and output sentence

example:

Input: "the quick brown fox",

>2. Neural network: structured with an equal number of input x and output y .

Note: in this example drawn a simple of uni-directional RNN, but in practice, this will usually be a bidirectional LSTM and bidirectional GRU and usually, a deeper model.

>3. Speech recognition property:

>>1. Number of input time steps here is very large ;

>>2. in speech recognition, usually the number of input time steps is much bigger than the number of output time steps.

e.g. have 10 seconds of audio and your features come at a 100 hertz so 100 samples per second, then a 10 second audio clip would end up with a 1000 inputs. But output might not have a thousand characters.

>4. CTC cost function: allows the RNN to generate an output like this repeated character: ttt, and special character "- blank character", and space.

e.g. output "ttt---h---eee---space---qqq-----". And, this is considered a correct output for the first parts of "the space q"

3. Speech recognition model

3.2: CTC cost model for speech recognitio

>5. Basic rule for the CTC cost function: to collapse (折叠) repeated characters not separated by "blank".

e.g. model output "ttt---h---eee---space---qqq-----" by collapsing repeated characters, not separated by blank, it actually collapse the sequence into "the space q".

and this allows network to have a thousand outputs by repeating characters a lot of times.

So, inserting a bunch of blank characters and still ends up with a much shorter output text transcript.

e.g. input sentence "the quick brown fox" including spaces actually has 19 characters, and if somehow, the newer network is forced to output a thousand characters by allowing the network to insert blanks and repeated characters , the output could still represent this 19 character with this 1000 outputs y

>6. Right bottom paper by Alex Grace, as well as by those deep speech recognition system, which I was involved in, used this idea to build effective Speech recognition systems.

(note: what is cost function for attention model /CTC, predict output is words' one-hot vector, labeled output is word one-hot vector? -use logistic regression for each element of one-hot vector? $I = \sum y_{i,j} \log(y^i_{j,i}) - (1-y_{i,j}) \log(1-y^i_{j,i})$ over i-word in vocab; or just care the label 1 position: $\sum y_{i,j} \log(y^i_{j,i})$ over i-word in vocab ??)-->should be softmax function output vocab size vector.

1.2 Training set

>3. Pre-process: audio clip:

Principle:

Even the human ear doesn't process raw wave forms, but the human ear has physical structures that measures the amounts of intensity of different frequencies, there is, a common pre-processing step for audio data .

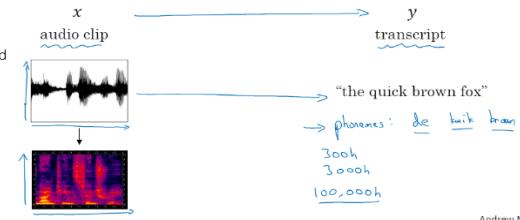
Pre-processing step for audio data :--false back outputs (伪空白输出)

Run raw audio clip and generate a spectrogram(声谱图) -right colored pic: the horizontal axis is time, and the vertical axis is frequencies, and intensity of different colors shows the amount of energy:

-->reflect: how loud is the sound at different frequencies and at different times

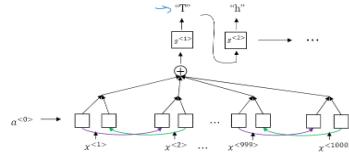
this process, also called false back outputs, is often commonly applied pre-processing step before audio is pass into the running algorithm. And the human ear does a computation very similar to this pre-processing step.

Speech recognition problem



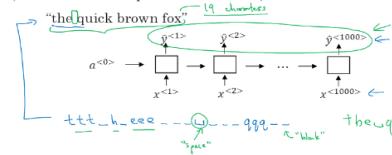
Andrew Ng

Attention model for speech recognition



CTC cost for speech recognition

(Connectionist temporal classification)



Basic rule: collapse repeated characters not separated by "blank" ↴

[Graves et al., 2006. Connectionist Temporal Classification: Labeling unsegmented sequence data with recurrent neural networks.] Andrew Ng

4. Summary:

This video shows a rough sense of how speech recognition models work. Attention like models work and CTC models work and present two different options of how to go about building these systems.

Now, today, building effective, production skills speech recognition system is a pretty significant effort and requires a very large data set. But, what I like to do in the next video is share you, how you can build a trigger word detection system, where keyword detection system which is actually much easier and can be done with even a smaller or more reasonable amount of data.

5.3-10: Trigger word detection_Audio data

You've now learned so much about deep learning and sequence models that we can actually describe a trigger word system quite simply just on one slide, as you see in this video.

1. Trigger word detection:

Definition:

trigger word detection systems: with the rise of speech recognition, there have been more and more devices you can wake up with your voice and those are sometimes called trigger word detection systems

Examples of trigger word systems:

Amazon Echo: woken up with the word "Alexa"
Baidu DuerOS: woken up with the phrase "xiaodunihao"
Apple Siri: working up with "hey, Siri"
Google Home: woken up with "okay, Google"

Application:

e.g.: an Amazon Echo in your living room, you can walk in your living room and just say, Alexa, what time is it.
e.g. 2: if you can build a trigger word detection system, maybe you can make your computer do something by telling it, computer, activate.
e.g. 3: One of my friends also worked on turning on and off a particular lamp using a trigger word kind of as a fun project.

2. Trigger word detection algorithm

The literature on trigger detection algorithm is still evolving, so there isn't wide consensus yet on what's the best algorithm for trigger word detection.

Examples of an algorithm can use:

2.1: Algorithm : setting target label to 1 only for trigger word (many to many architecture-Tx=Ty, standard RNN)

>1. Process:

>>>1. Input sentence pre-processing: computed spectrogram features of input audio clip: $x^{<1>} \dots x^{<Tx>}$
>>>2. Pass audio spectrogram feature $x^{<1>} \dots x^{<Tx>}$ through an RNN. (attention model or CTC)

>>>3. Define the target labels: 1 for trigger word, 0 for rest.

Idea: if some point in the input audio clip is when someone just finished saying the trigger word, then in the training sets:
>1. set the target labels to be 0 for everything before that point,
>2. right after that to set the target label of 1.
>3. And then if a little bit later on, the trigger word was said again at other time/ point, again set the target label to be 1 right after that.

>2. One slight disadvantage:

Model creates a very imbalanced training set to have a lot more 0 than 1. (consequences?) --> cost function: if output just all is 0, still get high precision.

(note: labeled y for one example: skewed: many 0, small 1: could not use accuracy to measure algorithm performance in dev, instead use F1, or others.)

Training example: one example:

input $x: x^{<1>} \dots x^{<Tx>}$ (there are negative words and positive words- very less positive words)
labeled output: $y: [0 \dots 0, 1, 1 \dots 0, 0 \dots 0, Ty]$ (for one positive word, output certain number like 50 of 1)
predict output y^{\wedge} (sigmoid): 0/1 for each $y^{<t>} \dots y^{<Ty>}$.

for one example, the prediction output is Ty : why only has like $50 \times$ positive word number of 1, result is 0-->accuracy on this example could not represent algorithm performance.

e.g. if just set prediction $y^{\wedge}=0 \dots$ if only have one positive word in example, then accuracy is $=Ty-50/Ty * 100$, still high (multi examples on dev set, -->average performance over all dev set examples)

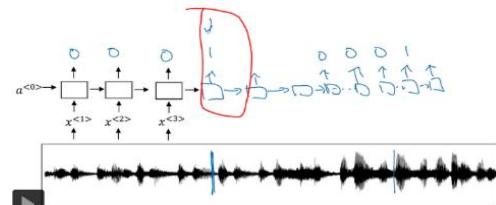
)

What is trigger word detection?



Trigger word detection algorithm

网易云课堂



Examples of an algorithm can use:

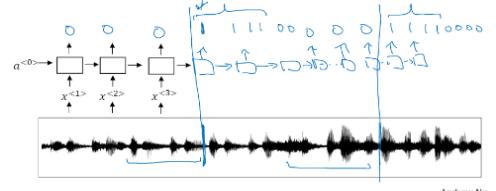
2.2: Algorithm : setting target label to multi 1 for trigger word

Idea:

Make above model a little bit easier to train: instead of setting only a single time step output 1, make it output a few 1 for several times or for a fixed period of time before reverting back to 0.

So that slightly evens out the ratio of 1 to 0, but this is a little bit of a hack.

Trigger word detection algorithm



Andrew Ng

3. Summarize:

In this course on sequence models, learned about our RNNs, including both GRUs and LSTMs, and then in the second week, learned a lot about word embeddings and how to learn representations of words. And then in this week, you learned about the attention model as well as how to use it to process audio data.

Summary

1. S2S model: many-many ($T_x = T_y$) RNN architecture: encoder + decoder application on machine translation , image captioning, etc.

2. Machine translation model: vs language model

>Language model: predict what's the next word given previous words: $p(y^{<1>} \dots y^{<t-1>})$
>>> standard language model: input $x^{<1>} \dots x^{<t-1>}$ = labeled $y^{<1>} \dots y^{<t-1>}$ (except $x^{<1>}=0, a^{<0>}=0$, vector)
>>> sampling standard large model: input $x^{<1>} \dots x^{<t-1>}$ = sampled, from prediction softmax distribution $y^{<t-1>}$.
>>> starting off activation $a^{<0>} = 0$

>Machine translation model: -->pick the maxim. translation $\max P(y|x)$

>>>architecture: many-many , encoder + decoder
>>> encoder: -->input sentence -->output representation of input sentence as starting off for decoder: context c
>>> decoder: similar with language model, only $a^{<0>} = 0$, but the representation of input, generated by encoder, -->also called conditional language model;
decoder input: = sampled previous step output $y^{<t-1>} +$ activation of previous step: $a^{<t-1>}$

3. Machine translation: choose best translation: $\max p(y|x)$

Method:

>1. greedy-->only pick max $p(y^{<1>} \dots y^{<t-1>} | x)$ at each time step -->could not guarantee get output y^{\wedge} is the max $p(y^{<1>} \dots y^{<t-1>} | x)$
>2. Beam search: could not guarantee find the best translation, but much better
Fine tuning beam search:
- Log: for numerical underflow;
- reduce penalty for long output: by divide Ty^{\wedge} / a for log p.

4. Error analysis for machine translation: RNN model (encoder + decoder) or beam search

Method:

- Take all the errors in dev/test set, where labeled output y^{\wedge} is provided by human.
- Get $p(y^{\wedge}|x)$ vs $p(y|x)$ through 'decoder'
- Compare: $p(y^{\wedge}|x)$ vs $p(y|x)$
- Summarize error fraction , caused by RNN and beam search-->get direction to work on

5. BLEU score: evaluation translation

Principle: check similarity /overlap with human translation
>basic precision for isolation word: for each output word, credit 1/0 if occur in ref.-->sum credit/ Ty
>modified precision for isolation word: p1:cout clip(ref)/cout (per se output)
> modified precision for n-gram word: pn
> Final bleu score: = $BP \times EXP(\text{average of } p1-pn)$
 BP : penalty for short output: =1 if output length > reference , otherwise <1: = exp (1 - reference length / output length)

6. Attention model: for very long sentence translation

S2S model, input all sentence (encoder) , then do the translation one word at a time. When input sentence is very long, hard to memorize enough info. of input sentence -->Bleu score reduce

Attention model: when output one word at a time, only look at certain input words
Architecture: (more context C to memorize input sentence info.-hidden layer info. of encoder) a little like deep model:

> first layer: BRNN: in each step, output feature $a^{<t>}$ of input word, instead of $y^{<t>}$; as input $x^{<1>} \dots x^{<t>}$ to second layer-context layer
> second layer-context layer: input: come from input word feature + output/decoder previous step activation $s^{<t-1>}$:
>>> assign weight to different input feature , get context C = sum , $a^{<t>} * a^{<t>}$ over all word 't' in input sentence
>>> weight $a^{<t>}$, use softmax to make sure sum over all input word t is 1 = $\exp e^{<t>}/\sum (\exp e^{<t>})$ all over input word
>>> $e^{<t,t>}$: function of output previous step status $s^{<t-1>}$, and feature of t'-th input word $a^{<t>}$ -use small neural network, train together in whole neural network.

>third layer-un-bidirectional RNN for output (only forward)
input: previous step output word $y^{<t-1>} + c^{<t>}$
output: activation $s^{<t>}$, output word $y^{<t>}$.

(note: context layer $c^{<1>} \dots c^{<t>}$ help to remember input sentence info., yet still need to read all input sentence, then go to next layer/step/before making prediction output)

7. Speech recognition:

output: is different from machine translation:

machine translation: output choose best during multi translations : max $P(y|x)$ -BLEU used to estimate output

speech recognition: recognize the correct character/word: output is right or not.

Property of speech recognition problem: $T_x \gg T_y$

model: 1. attention model

2. CTC model: output is character, input time step is very long

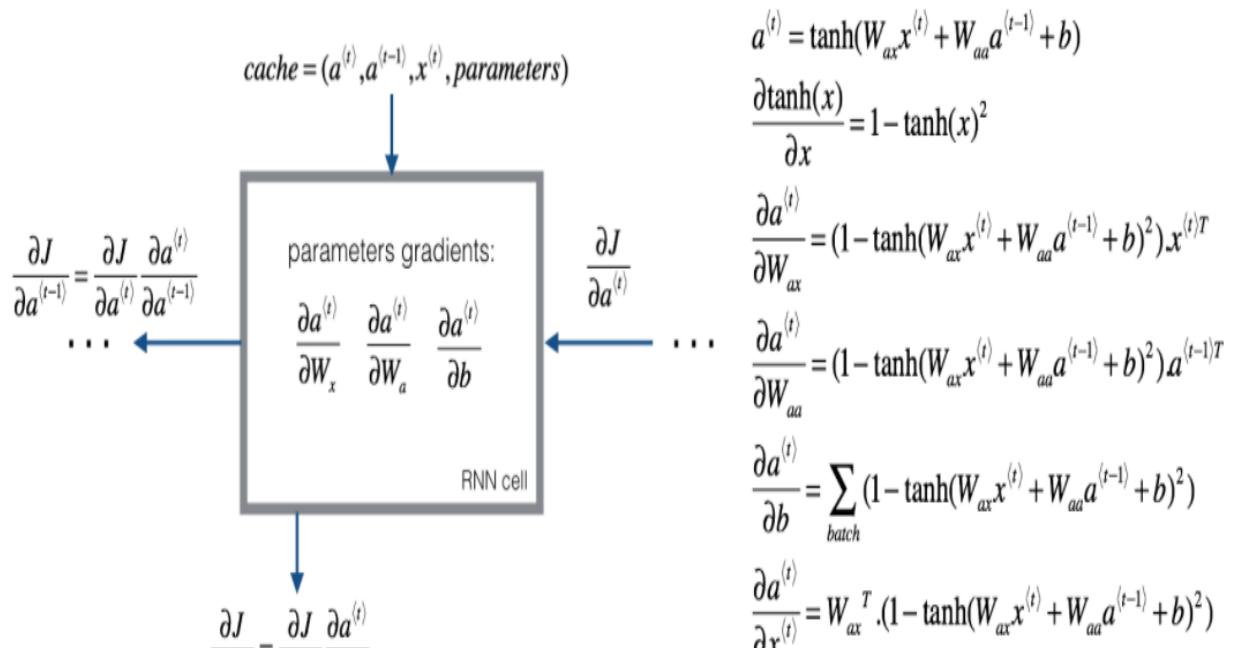
$T_x = T_y$, then collapse T_y

8. Trigger word: many-many sequence architecture:

model: output single 1 for trigger word/ output multi 1 for trigger word (even output 1 and 0)

3.1 - Basic RNN backward pass

We will start by computing the backward pass for the basic RNN-cell.



$$\frac{\partial \overline{x^{(t)}}}{\partial a^{(t)}} = \frac{\partial \overline{a^{(t)}}}{\partial x^{(t)}} \frac{\partial x^{(t)}}{\partial a^{(t)}}$$

$$\frac{\partial \overline{a^{(t)}}}{\partial a^{(t-1)}} = W_{aa}^T \cdot (1 - \tanh(W_{ax}x^{(t-1)} + W_{aa}a^{(t-1)} + b)^2)$$

http://blog.csdn.net/Koala_Tree

Figure 5: RNN-cell's backward pass. Just like in a fully-connected neural network, the derivative of the cost function J backpropagates through the RNN by following the chain-rule from calculus. The chain-rule is also used to calculate $(\frac{\partial J}{\partial W_{ax}}, \frac{\partial J}{\partial W_{aa}}, \frac{\partial J}{\partial b})$ to update the parameters (W_{ax}, W_{aa}, b_a) .

Deriving the one step backward functions:

To compute the `rnn_cell_backward` you need to compute the following equations. It is a good exercise to derive them by hand.

The derivative of \tanh is $1 - \tanh(x)^2$. You can find the complete proof [here](#). Note that: $\text{sech}(x)^2 = 1 - \tanh(x)^2$

Similarly for $\frac{\partial a^{(t)}}{\partial W_{ax}}$, $\frac{\partial a^{(t)}}{\partial W_{aa}}$, $\frac{\partial a^{(t)}}{\partial b}$, the derivative of $\tanh(u)$ is $(1 - \tanh(u)^2)du$.

The final two equations also follow same rule and are derived using the \tanh derivative. Note that the arrangement is done in a way to get the same dimensions to match.

Backward pass through the RNN

Computing the gradients of the cost with respect to $a^{(t)}$ at every time-step t is useful because it is what helps the gradient backpropagate to the previous RNN-cell. To do so, you need to iterate through all the time steps starting at the end, and at each step, you increment the overall db_a , dW_{aa} , dW_{ax} and you store dx .

Instructions:

Implement the `rnn_backward` function. Initialize the return variables with zeros first and then loop through all the time steps while calling the `rnn_cell_backward` at each time timestep, update the other variables accordingly.

```
# Increment global derivatives w.r.t parameters by adding their derivative at time-step t (~4 lines)
    dx[:, :, t] = dxt
    dwax += dwaxt
    dwaa += dwaat
    dba += dbat
```

3.2 - LSTM backward pass

3.2.1 One Step backward

The LSTM backward pass is slightly more complicated than the forward one. We have provided you with all the equations for the LSTM backward pass below. (If you enjoy calculus exercises feel free to try deriving these from scratch yourself.)

3.2.2 gate derivatives

$$d\Gamma_o^{(t)} = da_{next} * \tanh(c_{next}) * \Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}) \quad (7)$$

$$d\tilde{c}^{(t)} = dc_{next} * \Gamma_u^{(t)} + \Gamma_o^{(t)} (1 - \tanh(c_{next})^2) * i_t * da_{next} * \tilde{c}^{(t)} * (1 - \tanh(\tilde{c})^2) \quad (8)$$

$$d\Gamma_u^{(t)} = dc_{next} * \tilde{c}^{(t)} + \Gamma_u^{(t)} (1 - \tanh(c_{next})^2) * \tilde{c}^{(t)} * da_{next} * \Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}) \quad (9)$$

$$u_u - u_{next} + c_o \top o (\top - \text{tanh}(c_{next}) \top + c_o + u_{next} \top u \top (\top - \top u)) \quad (9)$$

$$d\Gamma_f^{(t)} = dc_{next} * \tilde{c}_{prev} + \Gamma_o^{(t)} (1 - \tanh(c_{next})^2) * c_{prev} * da_{next} * \Gamma_f^{(t)} * (1 - \Gamma_f^{(t)}) \quad (10)$$

3.2.3 parameter derivatives

$$dW_f = d\Gamma_f^{(t)} * \begin{pmatrix} a_{prev} \\ x_t \end{pmatrix}^T \quad (11)$$

$$dW_u = d\Gamma_u^{(t)} * \begin{pmatrix} a_{prev} \\ x_t \end{pmatrix}^T \quad (12)$$

$$\sqrt{\lambda} \nabla T$$

$$dW_c = d\tilde{c}^{(t)} * \begin{pmatrix} a_{prev} \\ x_t \end{pmatrix}^T \quad (13)$$

$$dW_o = d\Gamma_o^{(t)} * \begin{pmatrix} a_{prev} \\ x_t \end{pmatrix}^T \quad (14)$$

To calculate db_f, db_u, db_c, db_o you just need to sum across the horizontal (axis= 1) axis on $d\Gamma_f^{(t)}, d\Gamma_u^{(t)}, d\tilde{c}^{(t)}, d\Gamma_o^{(t)}$ respectively.

Note that you should have the `keep_dims = True` option.

Finally, you will compute the derivative with respect to the previous hidden state, previous memory state, and input.

$$da_{prev} = W_f^T * d\Gamma_f^{(t)} + W_u^T * d\Gamma_u^{(t)} + W_c^T * d\tilde{c}^{(t)} + W_o^T * d\Gamma_o^{(t)} \quad (15)$$

Here, the weights for equations 13 are the first n_a, (i.e. $W_f = W_f[:, n_a, :]$ etc...)

$$dc_{prev} = dc_{next} \Gamma_f^{(t)} + \Gamma_o^{(t)} * (1 - \tanh(c_{next})^2) * \Gamma_f^{(t)} * da_{next} \quad (16)$$

$$dx^{(t)} = W_f^T * d\Gamma_f^{(t)} + W_u^T * d\Gamma_u^{(t)} + W_c^T * d\tilde{c}_t + W_o^T * d\Gamma_o^{(t)} \quad (17)$$

where the weights for equation 15 are from n_a to the end, (i.e. $W_f = W_f[n_a :, :]$ etc...)

Exercise: Implement `lstm_cell_backward` by implementing equations 7 – 17 below. Good luck! :)

3.3 Backward pass through the LSTM RNN

This part is very similar to the `rnn_backward` function you implemented above. You will first create variables of the same dimension as your return variables. You will then iterate over all the time steps starting from the end and call the one step function you implemented for LSTM at each iteration. You will then update the parameters by summing them individually. Finally return a dictionary with the new

gradients.

Instructions: Implement the `lstm_backward` function. Create a for loop starting from T_x and going backward. For each step call `lstm_cell_backward` and update the your old gradients by adding the new gradients to them. Note that `dxt` is not updated but is stored.

```
# Loop back over the whole sequence
for t in reversed(range(T_x)):
    # Compute all gradients using lstm_cell_backward
    gradients = lstm_cell_backward(da[:, :, t] + da_prevt, dc_prevt, caches[t])
    # Store or add the gradient to the parameters' previous step's gradient
    dx[:, :, t] = gradients['dxt']
    dwf = dwf + gradients['dwf']
    dwi = dwi + gradients['dwi']
    dwc = dwc + gradients['dwc']
    dwo = dwo + gradients['dwo']
    dbf = dbf + gradients['dbf']
    dbi = dbi + gradients['dbi']
    dbc = dbc + gradients['dbc']
    dbo = dbo + gradients['dbo']
# Set the first activation's gradient to the uncomputed gradient da_prev.
```

```
# SET THE FIRST GRADIENT'S GRADIENT TO THE BACKPROPAGATED GRADIENT DA_PREV.  
da0 = gradients['da_prev']  
  
### END CODE HERE ###
```