Advanced Topics in Machine Learning

Alejo Nevado-Holgado

Lecture 12 (NLP 4) - Language models and vanilla RNNs V 0.3 (15 Feb 2020 - final version)

Course structure

- > Introduction: What is NLP. Why it is hard. Why NNs work well \leftarrow Lecture 9 (NLP 1)
- > Word representation: How to represent the meaning of individual words

 - Embeddings: First trick that boosted the performance of NNs in NLP Lecture 9 (NLP 1)
 - Word2vec: Single layer NN. CBOW and skip-gram ← Lecture 10 (NLP 2)
 - Co-occurrence matrices: Basic counts and SVD improvement ← Lecture 10 (NLP 2)
 - Glove: Combining word2vec and co-occurrence matrices idea ← Lecture 10 (NLP 2)
 - Evaluating performance of embeddings

 Lecture 10 (NLP 2)
- > Named Entity Recognition (NER): How to find words of specific meaning within text
 - Multilayer NNs: Margin loss. Forward- and back-propagation Lecture 11 (NLP 3)
 - Better loss functions: margin loss, regularisation \leftarrow Lecture 11 (NLP 3)
 - Better initializations: uniform, xavier ← Lecture 11 (NLP 3)
 - Better optimizers: Adagrad, RMSprop, Adam... ← Lecture 11 (NLP 3)

Course structure

> Language modelling: How to represent the meaning of full pieces of text

- Old technology: N-grams ← Lecture 12 (NLP 4)
- Recursive NNs language models (RNNs) ← Lecture 12 (NLP 4)
- Evaluating performance of language models \leftarrow Lecture 12 (NLP 4)
- Vanishing gradients: Problem. Gradient clipping Lecture 13 (NLP 5)
- Improved RNNs: LSTM, GRU ← Lecture 13 (NLP 5)
- > Machine translation: How to translate text
 - Old technology: Georgetown–IBM experiment and ALPAC report ← Lecture 16 (NLP 6)
 - Seq2seq: Greedy decoding, encoder-decoder, beam search \leftarrow Lecture 16 (NLP 6)

 - Evaluating performance: BLEU ← Lecture 16 (NLP 6)

The problem: Having a method to accurately represent the meaning of individual words (e.g. word embeddings), enormously helps to solve the simpler NLP tasks, but falls short in complex ones. Having a more advanced method that represents the meaning of full pieces of text (rather than separated words), helps to solve the more challenging tasks.

We created good word representations by predicting the presentation of the central word from the representations of a few neighbouring words (CBOW). We can create good representations of full pieces of text by predicting each word from all the preceding ones \rightarrow Language models



Language modelling

> A language model itself is simply an algorithm that assigns probabilities to the next



More formally: Given a sequence of words w⁽¹⁾, w⁽²⁾, w⁽³⁾, ..., w^(t) and their representations x⁽¹⁾, x⁽²⁾, x⁽³⁾, ..., x^(t), what the model calculates is an approximation of:

$$P(\boldsymbol{x}^{(t+1)} | \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(1)})$$

Language modelling

Applications: Predicting the next word in a piece of text has some direct applications



Applications: But predicting the next word in a piece of text is not the most useful thing itself. What is more useful is (1) the numerical representation that your method builds from all preceding text in order to accurately predict the next word, and (2) the method that accurately predicts the next word.

Both the presentation (1) and the method (2), when they solve accurately the language modelling problem, they also work well in other applications: NER, translation, question answering, conference resolution...



Old technology: N-grams

Idea: Assume that the probability of the next word, depends mostly from the previous 'n' words, rather than from the full preceding text. Then, rather than calculating the probability of every possible imaginable meaningful sequence of words, we only need to calculate the probability of each sequence of 'n' words

 $\mathsf{P}(\mathsf{x}^{(t+1)} \mid \mathsf{x}^{(t)}, \mathsf{x}^{(t-1)}, \mathsf{x}^{(t-2)}, ..., \mathsf{x}^{(1)}) \approx \mathsf{P}(\mathsf{x}^{(t+1)} \mid \mathsf{x}^{(t)}, \mathsf{x}^{(t-1)}, \mathsf{x}^{(t-2)}, ..., \mathsf{x}^{(t-n)})$

- Higher 'n' will theoretically be able to become more accurate... but would require very large amounts of data:
 - Unigrams: P("the"), P("students"), P("opened"), P("their")...
 - Bigrams: P("the students"), P("students opened"), P("opened their")...
 - Trigrams: P("the students opened"), P("students opened their")...
 - 4-grams: P("the students opened their")...

Old technology: N-grams

Implementation: For an n-gram of size 'n', count the number of times that each combination of 'n' and 'n-1' words appear in your corpus. Then apply the formula definition of "conditional probability"

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(x^{(t+1)}|x^{(t)}, x^{(t-1)}, \dots, x^{(t-n)}) = \frac{count(x^{(t+1)}, x^{(t)}, x^{(t-1)}, \dots, x^{(t-n)})}{count(x^{(t)}, x^{(t-1)}, \dots, x^{(t-n)})}$$

 $P(w \mid students opened their) = \frac{count(students opened their w)}{count(students opened their)}$

Old technology: N-grams

- > **Problems:** N-gram generated text is often surprisingly gramatical, but incoherent
 - **0** probabilities: P(w | students opened their) will be 0 if 'w' never appears in your corpus. You can partially solve this with smoothing, which adds a small number δ to all counts
 - Undefined probabilities: P(w | students opened their) will be '0/0' if 'students opened their' never appears in your corpus. You can partially solve this with backoff, using a 'n-1'-gram when the problem arises i.e. counting 'opened their'
 - Storage: You need to store counts for all 'n' and 'n-1' combinations. Size = $V^n + V^{n-1}$
 - Sparsity: Increasing 'n' to improve accuracy, makes the previous problems worse

 $P(w \mid students opened their) = \frac{count(students opened their w)}{count(students opened their)}$

10

Neural networks



nttps://cloud.google.com/blog/products/gcp/understandin g-neural-networks-with-tensorflow-playground 11

New technology: NNs

as the teacher started the clockthe students opened theirdiscard (out of window)process (window)

- Output distribution:
 - $\mathbf{y} = \operatorname{softmax}(\mathbf{U} \times \mathbf{h} + \mathbf{b}_2) \in \mathbb{R}^{\vee}$
- ➤ Hidden layer:

 $\mathbf{h} = f(\mathbf{W} \times \mathbf{e} + \mathbf{b}_1) \in \mathbb{R}^H$

Concatenated word embeddings:

 $\mathbf{e} = [\mathbf{e}^{(t-3)}, \mathbf{e}^{(t-2)}, \mathbf{e}^{(t-1)}, \mathbf{e}^{(t)}] = [\mathbf{e}^{(1)}, \mathbf{e}^{(2)}, \mathbf{e}^{(3)}, \mathbf{e}^{(4)}] \in \mathbb{R}^{4E}$

One-hot word vectors:

$$\mathbf{X} = [\mathbf{X}^{(t-3)}, \mathbf{X}^{(t-2)}, \mathbf{X}^{(t-1)}, \mathbf{X}^{(t)}] = [\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \mathbf{X}^{(3)}, \mathbf{X}^{(4)}] \in \mathbb{R}^{4\vee}$$



New technology: NNs

as the teacher started the clockthe students opened theirdiscard (out of window)process (window)

Advantages of window-NNs over n-grams:

- No sparsity problem
- Don't need to store all combinations of n words
- Remaining problems:
 - Fixed window width is too small... and we can never make it large enough!
 - Enlarging window increases the size of W
 - No symmetry in how words are processed: Each word (e.g. $x^{(1)}$ and $x^{(2)}$) are multiplied by completely different weights in W



➤ Idea: Recursively process one word at a time rather than all words within a window → No window limit, no increase in size of W, processing symmetry (all words are processed in the same manner)





Output distribution:

$$\mathbf{y}^{(t)} = \text{softmax}(\mathbf{U} \times \mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^V$$

➤ Hidden layer:

 $\mathbf{h^{(t)}} = f(\mathbf{W}_{\mathbf{h}} \times \mathbf{h^{(t-1)}} + \mathbf{W}_{\mathbf{e}} \times \mathbf{e^{(t)}} + \mathbf{b}_{\mathbf{1}}) \in \mathbb{R}^{H}$

- > One word embedding at a time: $e^{(t)} = E \times x^{(t)} \in \mathbb{R}^{E}$
- > One one-hot word vector at a time: $\mathbf{x}^{(t)} \in \mathbb{R}^{V}$



as the teacher started the clock the students opened their process

Advantages of RNNs over window-NNs:

- No window limit, can process any number of words
- No increase in size of W when processing more words
- Processing **symmetry**, all words are processed in the same manner
- Remaining problems:
 - Recurrent computation is slow
 - In practice, RNNs have problems using information from many words back (e.g. $x^{(t-100)}) \rightarrow$ short memory span



Learning in word2vec word model

- Word2vec very successfully implemented word embeddings using this context-meaning idea.
 - We start with a very large corpus of text (e.g. all of Wikipedia)
 - Every word is represented by a vector in \mathbb{R}^n space (n~200 dimentions)
 - You have a model (e.g. a NN) that tries to predict the vector of a word (i.e. the central word) given the vectors of the words around it (i.e. its context). In probability terms, the NN models the probability P(w_c | w_{c-3}, w_{c-2}, w_{c-1}, w_{c+1}, w_{c+2}, w_{c+3})
 - Go through each central word context pair in the corpus
 - In each iteration, modify the NN and vectors a little bit for words with similar contexts to have similar vectors
 - Repeat last 2 steps many times

Learning in RNNs language model

- Word2vec very successfully implemented word embeddings using this context-meaning idea.
 - We start with a very large corpus of text (e.g. all of Wikipedia)
 - Every word is represented by a vector in \mathbb{R}^n space (n~200 dimentions)
 - You have a model (e.g. a NN i.e. a RNN) that tries to predict the vector of a word (i.e. the central word i.e. the next word) given the vectors of the words around it (i.e. its context) before it. In probability terms, the NN models the probability $P(w_{c} + w_{c-3}, w_{c-2}, w_{c-1}, w_{c+1}, w_{c+2}, w_{t-1}, w_{t-2}, w_{t-3}, ...)$
 - Go through each central word context pair in the corpus
 - In each iteration, modify the NN and vectors a little bit for words with similar contexts to have similar vectors \rightarrow We now know we do this with backpropagation
 - Repeat last 2 steps many times

- > With backpropagation we want our RNN to gradually learn to approximate the probability $P(w_{t+1} | w_{t}, w_{t-1}, w_{t-2}, w_{t-3}, ...)$
- > But backpropagation needs an error function that represents numerically how far our RNN is from perfectly approximating P($w_{t+1} | w_{t}, w_{t-1}, w_{t-2}, w_{t-3}, ...$)
- Error functions used before:
- CBOW (word2vec): Negative Sampling^[arXiv:1310.4546] (simplified version of Negative Contrastive Sampling^[arXiv:1206.6426])

$$J = \log sig(\underbrace{[y_{centre}]}_{E} \cdot \underbrace{[e_{context}]}_{E}) + \sum_{i \sim P_{noise}}^{K} \log sig(\underbrace{[y_i]}_{E} \cdot \underbrace{[e_{context}]}_{E}))$$

 $[something] \longrightarrow `something' is a vector of V elements$ $[something](v) \longrightarrow element `v' of vector [something]$ $V V [y_{something}] - [y_{so$

 $\begin{array}{c} E \longrightarrow \text{number of embedding dimensions} \\ [e_{something}] \longrightarrow \text{embedding of 'something'} \\ [y_{something}] \longrightarrow \text{output of CBOW NN when input is } [e_{something}] \\ E \end{array}$

- More error functions used before:
- Glove (hybrid of NN and co-occurrence matrix): Error function based on 3 heuristic considerations

$$J = \sum_{i,j=1}^{V} ramp\left(\begin{bmatrix} C \\ V, V \end{bmatrix} (i,j) \right) \left(\begin{bmatrix} e_i \\ E \end{bmatrix} \cdot \begin{bmatrix} e_j \end{bmatrix} + b_i + b_j - log\left(\begin{bmatrix} C \\ V, V \end{bmatrix} (i,j) \right) \right)^2$$

 $V \longrightarrow$ size of vocabulary $b_i \longrightarrow$ arbitrary constant for word 'i' $ramp(\dots) \longrightarrow$ almost a ramp function

 $\begin{bmatrix} C \\ V, V \end{bmatrix}$ —> co-occurrence count matrix $\begin{bmatrix} C \\ V, V \end{bmatrix}$ (*i*, *j*) \longrightarrow count of how many times word 'i' is near word 'j'

- NER: Margin loss

$$J = max(0, 1 + \underline{s_{NEG}} - \underline{s_{POS}})$$

 $s_{NEG} \longrightarrow$ output of NER NN when input is $\begin{bmatrix} x_{NEG} \\ 5E \end{bmatrix}$ $s_{POS} \longrightarrow$ output of NER NN when input is $\begin{bmatrix} x_{POS} \\ 5E \end{bmatrix}$

- > New error function:
- Cross-entropy

$$J = -\begin{bmatrix} x_{t+1} \end{bmatrix} \cdot log(\begin{bmatrix} y_t \end{bmatrix}) = -log(\begin{bmatrix} y_t \end{bmatrix}(i))$$









25



- ➤ Minibatch learning: Computing ∑J^(t)(θ) over all possible values of t (i.e. over all word positions in the corpus) in each iteration of Stochastic Gradient Descent (SGD) is too expensive. We rather computer ∑J^(t)(θ) only for a few random sentences in each iteration of SGD.
- **Backpropagation through time:** If we unfold the RNN along the previous time steps (t-1, t-2, t-3, ...), we can apply the chain rule to calculate the gradient of $\sum J^{(t)}(\theta)$ like in any other standard NN



The chain rule in RNNs: Given a function f(x,y) that is a function of other two different functions of 't' x(t) and y(t) (making f(x,y) = f(x(t),y(t))), according to the chain rule the gradient of f is:

$$\underbrace{\frac{d}{dt}f(\boldsymbol{x}(t),\boldsymbol{y}(t))}_{dt} = \frac{\partial f}{\partial \boldsymbol{x}}\frac{d\boldsymbol{x}}{dt} + \frac{\partial f}{\partial \boldsymbol{y}}\frac{d\boldsymbol{y}}{dt}$$

Derivative of composition function



The chain rule in RNNs: Given a function f(x,y) that is a function of other two different functions of 't' x(t) and y(t) (making f(x,y) = f(x(t),y(t))), according to the chain rule the gradient of f is:

$$\underbrace{\frac{d}{dt}f(\boldsymbol{x}(t),\boldsymbol{y}(t))}_{\underline{dt}} = \frac{\partial f}{\partial \boldsymbol{x}}\frac{d\boldsymbol{x}}{dt} + \frac{\partial f}{\partial \boldsymbol{y}}\frac{d\boldsymbol{y}}{dt}$$

Derivative of composition function



The chain rule in RNNs: Given a function f(x,y) that is a function of other two different functions of 't' x(t) and y(t) (making f(x,y) = f(x(t),y(t))), according to the chain rule the gradient of f is:

$$\underbrace{\frac{d}{dt}f(\boldsymbol{x}(t),\boldsymbol{y}(t))}_{\underline{dt}} = \frac{\partial f}{\partial \boldsymbol{x}}\frac{d\boldsymbol{x}}{dt} + \frac{\partial f}{\partial \boldsymbol{y}}\frac{d\boldsymbol{y}}{dt}$$

Derivative of composition function



Once trained on a large corpus, you can generate text with a RNN language model by sampling the output distribution y^(t) in each time step t, and using that word for time step t+1



> RNN trained on obama speeches^[Samin @ Medium]:

The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done.

➢ RNN trained on Harry Potter^[Max Deutsch @ Medium]:

"Sorry," Harry shouted, panicking—"I'll leave those brooms in London, are they?"

"No idea," said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry's shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn't felt it seemed. He reached the teams too.





> RNN trained on paint color names^[Janelle Shane]:

Ghasty Pink 231 137 165 Power Gray 151 124 112 Navel Tan 199 173 140 Bock Coe White 221 215 236 Horble Gray 178 181 196 Homestar Brown 133 104 85 Snader Brown 144 106 74 Golder Craam 237 217 177 Hurky White 232 223 215 Burf Pink 223 173 179 Rose Hork 230 215 198

Sand Dan 201 172 143
Grade Bat 48 94 83
Light Of Blast 175 150 147
Grass Bat 176 99 108
Sindis Poop 204 205 194
Dope 219 209 179
Testing 156 101 106
Stoner Blue 152 165 159
Burble Simp 226 181 132
Stanky Bean 197 162 171
 Turdly 190 164 116

> This RNNs runs at the character level, not word level

Evaluating quality of language models

Perplexity: This is the standard evaluation metric for language models. It measures how 'surprised' the NN is of finding that the real word at step t+1 is x^(t+1) after reading as input all previous words x^(t), x^(t-1), x^(t-2), x^(t-3), ...

$$perplexity(t) = \prod_{t=1}^{T} \left(\frac{1}{P_{RNN} \left(x^{(t+1)} \mid x^{t}, x^{t-1}, x^{t-2}, \dots \right)} \right)^{1/T} = \prod_{t=1}^{T} \left(\frac{1}{\begin{bmatrix} x_{t+1} \end{bmatrix} \cdot \begin{bmatrix} y_{t} \end{bmatrix}} \right)^{1/T}$$

Surprisingly enough, if we apply the exponential rules e^{ab}=e^a+e^b and e^{log(a)}=a, it turns out that perplexity is equal to the exponential of the total loss J:

$$\prod_{t=1}^{T} \left(\frac{1}{\begin{bmatrix} \boldsymbol{x_{t+1}} \end{bmatrix} \cdot \begin{bmatrix} \boldsymbol{y_t} \end{bmatrix}} \right)^{1/T} = exp \left(\frac{1}{T} \sum_{t=1}^{T} -log \left(\begin{bmatrix} \boldsymbol{x_{t+1}} \end{bmatrix} \cdot \begin{bmatrix} \boldsymbol{y_t} \end{bmatrix} \right) \right) = exp \left(J(\theta) \right)$$

Evaluating quality of language models

> Language models have greatly improved perplexity during the last years^[Grave & Joulin]

Model	Perplexity
Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
gly Sparse Non-negative Matrix factorization (Shazeer et al., 2015)	52.9
LSTM-2048 (Jozefowicz et al., 2016)	43.7
2-layer LSTM-8192 (Jozefowicz et al., 2016)	30
Ours small (LSTM-2048)	43.9
Ours large (2-layer LSTM-2048)	39.8

(lower is better)

[Gave & Joulin] = https://research.fb.com/building-an-efficient-neural-language-model-over-a-billion-words/

- > Language models become a subcomponent task of many other more complex tasks:
 - Predictive typing
 - Speech recognition
 - Handwriting recognition
 - Spelling/grammar correction
 - Text classification
 - Machine translation
 - Summarization
 - Dialogue
 - etc...

Language models can be used for Part Of Speech (POS) or Named Entity Recognition (NER):



> Language models can be used for text classification (e.g. sentiment analysis):



Language models can be used for question answering



Next lecture!

- > In this lecture: We have described the simplest type of RNN \rightarrow "vanilla RNN"
- > Next lecture: We will learn more complex and powerful ones



 \rightarrow Multiplayer RNNs



➢ By end of course: We will learn much more advanced ones

bidirectional attention self-attention residual networks transformers reformers...



Course structure

- > Introduction: What is NLP. Why it is hard. Why NNs work well \leftarrow Lecture 9 (NLP 1)
- > Word representation: How to represent the meaning of individual words
 - Old technology: One-hot representations, synsets ← Lecture 9 (NLP 1)
 - Embeddings: First trick that boosted the performance of NNs in NLP ← Lecture 9 (NLP 1)
 - Word2vec: Single layer NN. CBOW and skip-gram ← Lecture 10 (NLP 2)
 - Co-occurrence matrices: Basic counts and SVD improvement ← Lecture 10 (NLP 2)
 - Glove: Combining word2vec and co-occurrence matrices idea ← Lecture 10 (NLP 2)
 - Evaluating performance of embeddings

 Lecture 10 (NLP 2)
- > Named Entity Recognition (NER): How to find words of specific meaning within text
 - Multilayer NNs: Margin loss. Forward- and back-propagation Lecture 11 (NLP 3)
 - Better loss functions: margin loss, regularisation ← Lecture 11 (NLP 3)
 - Better initializations: uniform, xavier ← Lecture 11 (NLP 3)
 - Better optimizers: Adagrad, RMSprop, Adam... ← Lecture 11 (NLP 3)

Course structure

> Language modelling: How to represent the meaning of full pieces of text

- Old technology: N-grams ← Lecture 12 (NLP 4)
- Recursive NNs language models (RNNs) ← Lecture 12 (NLP 4)
- Evaluating performance of language models \leftarrow Lecture 12 (NLP 4)
- Vanishing gradients: Problem. Gradient clipping Lecture 13 (NLP 5)
- Improved RNNs: LSTM, GRU ← Lecture 13 (NLP 5)
- > Machine translation: How to translate text
 - Old technology: Georgetown–IBM experiment and ALPAC report ← Lecture 16 (NLP 6)
 - Seq2seq: Greedy decoding, encoder-decoder, beam search \leftarrow Lecture 16 (NLP 6)

Literature

- > Papers =
 - "N-gram language models", Jurafsky et al., 2018.
 <u>https://web.stanford.edu/jurafsky/slp3/3.pdf</u>
 - "The unreasonable effectiveness of recurrent neural networks", Karpathy, 2015. http://karpathy.github.io/2015/05/21/rnn-effectiveness/
 - Sections 10.1 and 10.2 of "Deep Learning", Goodfellow et al., 2016. http://www.deeplearningbook.org/contents/rnn.html