NEURAL MODELS FOR WORD EMBEDDING

Prof. Eduardo Bezerra (CEFET/RJ) ebezerra@cefet-rj.br





- Neural Nets (quick overview)
 Word2vec
 - Language models
 - Skip-gram



Model of a biological neuron



http://creationwiki.org/Neuron

Artificial Neuron

- 5
- □ The artificial neuron is a model **inspired** by the real
- □ In an artificial neuron:
 - input are features
 - each feature has an associated weight
 - weighted sum of features is called neuron pre-activation
 - value produced in pre-activation goes through a nonlinearity to produce neuron activation

Artificial Neuron

Ь (x_0) The artificial neuron is a model inspired by the real one (biological neuron). w_1 (x_1) w_2 $a = \hat{y}$ $z = w^T x + b \quad a = \varphi(z)$ (x_2) w_3 (x_3) Wn (x_n)

Artificial Neuron - input



Artificial Neuron – weights



Artificial Neuron – pre-activation



pre-activation

Artificial Neuron – activation function



activation

Nonlinearities are required so that the network can learn complex representations of the data.

> "The composition of linear transformations is also a linear transformation"

Artificial Neuron – activation functions



4

2

Z,



ReLU(z)

(a) Logística.

-4 - 2 0

(b) Hiperbólica tangente.

-4 - 2 0

Ζ.

2 4

(c) Retificada linear.

Artificial Neuron – activation functions

- Sigmoid Logistic
- Tangent Hyperbolic
- ReLU
- softmax

$\varphi(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$
$\varphi(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$
$\varphi(z) = \operatorname{ReLU}(z) = \max(0, z)$
$y_i = o\left(a_i^{(L+1)}\right) = \frac{e^{a_i^{(L+1)}}}{\sum_{j=1}^C e^{a_j^{(L+1)}}}$



$$y_i = o\left(a_i^{(L+1)}\right) = \frac{e^{a_i^{(L+1)}}}{\sum_{j=1}^C e^{a_j^{(L+1)}}}$$

13

Maps the NN outputs to a probability distribution.





softmax

$$y_i = o\left(a_i^{(L+1)}
ight) = rac{e^{a_i^{(L+1)}}}{\sum_{j=1}^C e^{a_j^{(L+1)}}}$$

Maps the NN outputs to a probability distribution.



Fonte: https://www.tensorflow.org/versions/r0.9/tutorials/mnist/beginners/index.html

Artificial Neural Net

It is possible to architect arbitrarily complex networks using the artificial neuron as the basic component.





Artificial Neural Net





Feedforward Neural Network

. . .

. . .

. . .

Kinds of layers

□ In an ANN:

input layer are the input patterns (features);

output layer is result of the computation performed;

other layers are called hidden layers.

Amount of layers and amount of units (neurons) per layer are part of the ANN architecture.

Kinds of layers - example



Kinds of ANNs

Fully connected x not fully connected

Recurrent x feed forward





ANNs - notation

Pre-activation:

- $\begin{array}{ll} z_i^{[l]} & \leftarrow \text{identifica a camada} \\ \leftarrow \text{identifica o neurônio dentro da camada} \ l \end{array}$

Activation:

- $\begin{array}{ll} a_i^{[l]} & \leftarrow \text{identifica a camada} \\ \leftarrow \text{identifica o neurônio dentro da camada} \ l \end{array}$

ANNs – notation - example

e.g., in a two-layer ANN:



Computations in the first layer:

$$\begin{aligned} z_1^{[1]} &= (w_1^{[1]})^T x + b_1^{[1]}, a_1^{(1)} = \varphi(z_1^{[1]}) \\ z_2^{[1]} &= (w_2^{[1]})^T x + b_2^{[1]}, a_2^{(1)} = \varphi(z_2^{[1]}) \\ z_3^{[1]} &= (w_3^{[1]})^T x + b_3^{[1]}, a_3^{(1)} = \varphi(z_3^{[1]}) \\ z_4^{[1]} &= (w_4^{[1]})^T x + b_4^{[1]}, a_4^{(1)} = \varphi(z_4^{[1]}) \end{aligned}$$

Rectified linear units, ReLU

- Convnet of 4 layers with ReLUs.
 - CIFAR-10 (60K, 32x32 imgs)
 - It reaches the same error level 6 times faster than tanh.



ImageNet Classification with Deep Convolutional Neural Networks, 2012.



Cost function

- 24
- The error signal is the difference between y⁽ⁱ⁾ and output produced for x⁽ⁱ⁾.
 - measures the difference between the <u>network predictions</u> and the <u>desired values</u>.
- A cost function must provide a way to measure the error signal for a set o training examples.

Cost function J(W,b)

- Cost function: error of the network considering all the examples.
 - function of the set of all weights and bias of the network.



Some (unregularized) cost functions

26

$$\begin{aligned} \mathcal{D} &= \{ (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) : 1 \leq i \leq m \} \\ J_{\text{MSE}} &= \frac{1}{m} \sum_{i=1}^{m} \left[f(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)} \right]^2 & \text{Mean square error} \\ J_{\text{cross entropy}} &= \sum_{i=1}^{m} \left[p(\mathbf{x}^{(i)}) \log q(\mathbf{x}^{(i)}) \right] & \text{Cross entropy error} \\ J_{\text{KL}} &= \sum_{i=1}^{m} \left[p(\mathbf{x}^{(i)}) \log \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} \right] & \text{Kullback-Leibler divergence} \end{aligned}$$



Training

□ Given a training set of the form

 $\{\mathbf{x}^{(i)}, \mathbf{y}^{(i)} : 1 \le i \le m\}$

training an ANN corresponds to using this set to adjust the parameters of the network, so that the training error is minimized.

□ So, training of an ANN is an optimization problem.

Training

29

The error signal (computed with a cost function) is used during training to gradually change the weights (parameters), so that the predictions are more accurate.



Stochastic Gradient Descent (SGD)

30

SGD allows to traverse the space of parameters of a neural network in search for a local minimum.





Source: David McKay, ITPRNN

Alternatives to SGD



Source: http://sebastianruder.com/optimizing-gradient-descent/

³² Neural Nets - Backpropagation

Error back-propagation

- Consider that the output produced for x⁽ⁱ⁾ is different from y⁽ⁱ⁾ (i.e., there is an nonzero error signal).
- Credit Assignment Problem
 - Then it is necessary to determine the responsibility of each parameter (weight) of the network by this error...
 - ... and change these parameters for the purpose of reducing the error.

Error back-propagation

In an ANN with at least one hidden layer, there is no direct supervision signal for these layers.



Error back-propagation

35

This problem gets worse as the amount of hidden layers increases.





Backpropagation

- 36
- Algorithm invented many times...
- Recursively propagates the error signal from the output layer through the hidden layers, to the input layer.
- Used in conjunction with some optimization algorithm (e.g., SGD) to gradually minimize network error.
Backpropagation

Backpropagation, an abbreviation for "backward propagation of errors", is a common method of training artificial neural networks used in conjunction with an optimization method such as gradient descent. The method calculates the gradient of a loss function with respect to all the weights in the network. The gradient is fed to the optimization method which in turn uses it to update the weights, in an attempt to minimize the loss function. --Frederik Kratzert

ANN training with Backpropagation

- 38
- Set initial weights for the network
- Present an example and get the result
- □ For each layer (from the last):
 - Calculate the error for each neuron
 - Update the weights (with SGD or variant)
 - Propagate this error to the previous layer
- Repeat with other examples until convergence

Backpropagation – example



Fonte: Automatic Differentiation in Machine Learning: a Survey (https://arxiv.org/pdf/1502.05767.pdf)





Credit where credit is due

- □ The original papers from Mikolov et al.
- CS 224D: Deep Learning for NLP
 https://cs224d.stanford.edu/lecture_notes/notes1.pdf
- □ McCormick, C. (2016). Word2Vec Tutorial.
 - <u>http://www.mccormickml.com</u>
- Understanding word vectors (Python notebook)
 Inderstanding word vectors (Python notebook)
 - https://gist.github.com/aparrish/2f562e3737544cf29aaf1 af30362f469

- A model that assigns a probability to a sequence of tokens.
- □ A good language model gives...
 - ...(syntactically and semantically) valid sentences a high probability.
 - ...low probability to nonsense.

Mathematically, we can apply a LM to any given sequence of n words:

$$P(w_1, w_2, \cdots, w_n)$$

- 45
- □ An example:
 - "The quick brown fox jumps over the lazy dog."
- Another example:

"The <u>quik</u> brown <u>lettuce</u> <u>over jumps</u> the lazy dog."



46

$$P(w_1,w_2,\cdots,w_n)=\prod_{i=1}^n P(w_i)$$
 Unigram model

$$P(w_1,w_2,\cdots,w_n)=\prod_{i=2}^n P(w_i|w_{i-1})$$
 Bigram mode

But, how to learn these probabilities?





https://research.fb.com/people/mikolov-tomas/

- Efficient Estimation of Word Representations in Vector Space, September 7th, 2013.
- Distributed Representations of Words and Phrases and their Compositionality, October 16th, 2013.

https://code.google.com/archive/p/word2vec/

The word2vec trick

- word2vec uses a trick:
 - A single hidden layer neural network is trained to perform a certain "fake" task.
 - But this NN is not actually used!
- Instead, the goal is to learn the weights of the hidden layer- these weights are the "word vectors".

The fake tasks

Skip-gram: predicting surrounding <u>context words</u> given a <u>center word</u>.

CBOW: predicting a center word from the surrounding context.



Skip-gram

- - The task: given a specific word w in the middle of a sentence (the input word), look at the words <u>nearby</u> and pick one word at random.
 - We then train the network to produce the probability (for every word in the vocabulary) of being <u>nearby</u> w.
 - "nearby" means there is actually a "window size" hyperparameter (typical value: 5)

Skip-gram

- 52
- Output probabilities are going to relate to how likely it is to find each vocabulary word nearby our input word w.
- □ For example, say **w** = "Africa"
 - output probabilities should be much higher for related words (e.g. "lion", "zebra") than for unrelated words (e.g. "bear", "kangaroo").

Training examples (word pairs)

- The training examples will be word pairs taken from the input corpus.
- During training, the network will learn parameters that capture the statistics related to the number of times each pairing shows up.

Training examples (word pairs)

- For example, suppose the network gets many more training samples of ("Soviet", "Union") than ("Soviet", "Sasquatch").
 - After training, if the NN is given the word "Soviet", it will output a much higher probability for "Union" (or "Russia") than it will for "Sasquatch".

Training examples (word pairs)

Source Text	Training Samples
The quick brown fox jumps over the lazy dog. \implies	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. \Longrightarrow	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. \Longrightarrow	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. \implies	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

"The quick brown fox jumps over the lazy dog." window size = 2.

The word highlighted in blue is the input word **w**.

But, during training, we cannot provide the NN with the text of the words.

□ We need to use some <u>vector representation</u> for each word in the vocabulary → so, we have to use one hot encoding.

Rome = [1, 0, 0, 0, 0, 0, 0, ..., 0]Paris = [0, 1, 0, 0, 0, 0, ..., 0]Italy = [0, 0, 1, 0, 0, 0, ..., 0]France = [0, 0, 0, 1, 0, 0, ..., 0]





Credits: Marco Bonzaninin

Skip-gram NN architecture



"Linear neurons" means there is no activation function on the hidden layer neurons, ...

Skip-gram NN architecture



..., but the output neurons use softmax.

Skip-gram NN architecture



The amount of neurons in the hidden layer (a hyperparameter) determines de size of the embedding.

62

The hidden layer

Multiplying a 1 x n one-hot vector v by an n x d matrix W (the hidden layer) will effectively select the row of W corresponding to the "1" position in v.

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = \begin{bmatrix} 10 & 12 & 19 \end{bmatrix}$$

So, the output of the hidden layer is just the "word vector" for the input word \rightarrow the hidden layer plays the role of a **lookup table**.

The hidden layer

64



The output layer

65

The 1 x d word vector (coming from the hidden layer) then gets fed to the output layer.
 The output layer has |V| neurons and is a softmax regression classifier.



Output weights for "car"

Probability that if you

 randomly pick a word nearby "ants", that it is "car"

Skip-gram captures context similarity

- If two words w_i and w_k have similar "contexts" (i.e., similar words are likely to appear around them), then the model needs to output very similar results for them.
 - One way for the network to do this is if the word vectors for w_i and w_k are similar.
- So, if two words have similar contexts, our network is motivated to learn similar word vectors for them.

A drawback of skip-gram

- The resulting trained model contains a huge amount of weights.
 - e.g., d = 300 and |V| = 10,000, that's 3M weights in the hidden layer and output layer each!
- Training (optimizing) a skip-gram model on a large corpus would be prohibitive.
- □ In order to circumvent this, some innovations...

Three inovations

- 1. Learn phrases from the corpus.
- 2. Subsample frequent words to decrease the number of training examples.
- 3. Modify the optimization objective with a technique called "Negative Sampling", which causes each training sample to update only a small percentage of the model's weights.

Learning Phrases

- A word pair like "Boston Globe" (a newspaper) has a much different meaning than the individual words "Boston" and "Globe".
- So it makes sense to treat any occurrences of "Boston Globe" as a single word with its own word vector representation.

Learning Phrases

- \Box In an experiment from the 2nd word2vec paper:
 - a model was trained on <u>100 billion words</u> from a Google News dataset, using d=300.
 - the addition of phrases to the model reduced the vocabulary size to <u>3 million "words"</u>.
- □ This (pre-trained) model can be freely downloaded.
- Phrase detection is also covered in the paper.

Learning Phrases

- Each pass of their phrase detection tool only looks at combinations of 2 words.
- But it can be executed multiple times to get longer phrases.
- So, the first pass will pick up the phrase "New_York", and then running it again will pick up "New_York_City" as a combination of "New_York" and "City".

Subsampling frequent words

- Two issues with common words like "the":
 - When looking at word pairs, ("fox", "the") doesn't tell us much about the meaning of "fox".
 - 2. We will have many more samples of ("the", ...) than we need to learn a good vector for "the".
- □ A "subsampling" scheme addresses this...
Subsampling of frequent words

- Example: for a window size of 10, if a specific instance of "the" is removed:
 - 1. As we train on the remaining words, "the" will not appear in any of their context windows.
 - 2. We'll have 10 fewer training samples where "the" is the input word.
- □ This has a huge effect on the size of the training set.

Subsampling of frequent words

- "subsampling" scheme:
 - For each word found in the corpus, there is a chance that it will be discarded.
 - The probability that a word w_i will be discarded is related to its frequency:

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$

Subsampling of frequent words





- Train a NN means to iterate on: "take each x⁽ⁱ⁾ and adjust <u>all of the weights</u> slightly so that the NN predicts x⁽ⁱ⁾ more accurately".
 - i.e., each training sample will tweak all of the weights in the neural network.
 - But the word2vec architecture is fully connected!
- Negative sampling addresses this problem.

- Basic idea: modify only a <u>small percentage</u> of the weights in the NN for each training sample.
- When training the network on the word pair ("fox", "quick"), recall that the "label" (response signal) is a one-hot vector.

- With negative sampling, just a small number of <u>"negative" words</u> (let's say 5) are randomly selected to update the weights.
 - (a "negative" word is one for which we want the network to output a 0 for).
- We still update the weights for the "positive" word.

Negative sampling - example

- Consider Negative Sampling applied to a output layer with a 300 x 10,000 weight matrix.
 - we will just be updating the weights for the positive word, plus the weights for 5 other words that we want to output 0.
 - That's a total of 6 output neurons, and 1,800 weight values.
 That's only 0.06% of the 3M weights in the output layer!

The "negative samples" are chosen using a "unigram distribution".

The probability for selecting a word as a negative sample is related to its frequency, with more frequent words being more likely to be selected as negative samples.

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=0}^n \left(f(w_j)^{3/4}\right)}$$

□ The 2nd paper says that

- selecting 5-20 words works well for smaller datasets,
- only 2-5 words suffice for large datasets.