

CS657A: INFORMATION RETRIEVAL NEURAL NETWORK BASED EMBEDDING MODELS

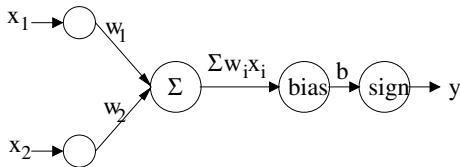
Arnab Bhattacharya
`arnabb@cse.iitk.ac.in`

Computer Science and Engineering,
Indian Institute of Technology, Kanpur
<http://web.cse.iitk.ac.in/~cs657/>

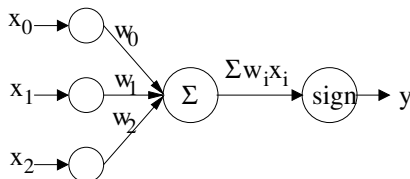
2nd semester, 2021-22
Tue 1030-1145, Thu 1200-1315

Perceptron

- A **perceptron** is a simple binary *linear* classifier
- Input attributes x_1, \dots, x_n are weighted and summed
- A bias b is added as well
- Final class ($y = \pm 1$) is *sign* of the output
- The *sign* node is called the **activation function** or **link/decision/transfer function**
- Decision boundary is $\vec{w} \cdot \vec{x} + b = \sum_{i=1}^n w_i x_i + b$
- Therefore, **sign** of $\vec{w} \cdot \vec{x} + b$ predicts the class



- Why is bias needed?
- Otherwise, hyperplane passes through origin
- Simple trick to model input uniformly: include 1 as x_0 of data
- Decision boundary becomes $w_0x_0 + \vec{w} \cdot \vec{x} = \sum_{i=0}^n w_i x_i$
- Weight on $x_0 = 1$ becomes the constant term, i.e., $w_0 = b$



Examples of Perceptrons

- Different boolean functions
- AND (of x_1 and x_2)
 - $w_1 = w_2 = 1, w_0 = -1.5$
- OR (of x_1 and x_2)
 - $w_1 = w_2 = 1, w_0 = -0.5$
- NOT (of x_1)
 - $w_1 = -1, w_0 = 0.5$
- XOR (of x_1 and x_2)
 - Cannot be done as the two classes are not linearly separable

Learning a Perceptron

- What is new in a linear classifier?
- Training a perceptron, i.e., learning the weights w
- **Perceptron learning rule** or **perceptron training rule**

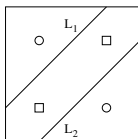
$$w_i = w_i - \eta(\hat{y}_i - y_i)x_i$$

where \hat{y}_i is the predicted value and η is the *learning rate*

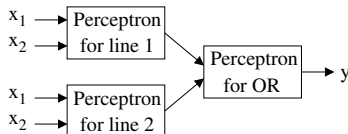
- If $y_i = \hat{y}_i$, there is no change in weight
- If $y_i = +1$ and $\hat{y}_i = -1$, i.e., $\hat{y}_i - y_i < 0$
 - Weights of positive x_i are increased and those of negative x_i are decreased thereby pushing \hat{y}_i towards positive
- If $y_i = -1$ and $\hat{y}_i = +1$, i.e., $\hat{y}_i - y_i > 0$
 - Weights of positive x_i are decreased and those of negative x_i are increased thereby pushing \hat{y}_i towards negative
- If the data *is* linearly separable, a perceptron *will* learn it, i.e., it will converge to the global optimum; otherwise, it may oscillate
- The learning rates η may be modified to give more importance to recent examples
- Weights can be learned through *gradient descent* method as well

Combination of Perceptrons

- A single perceptron can learn only a single hyperplane
- If the data can be separated using two or more hyperplanes, a *combination* of perceptrons can learn it
- Example: XOR

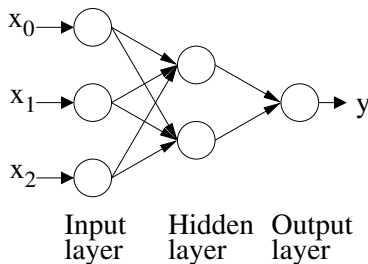


- Each hyperplane can be modeled by a perceptron and the outputs can be combined using OR



Artificial Neural Networks

- **Artificial neural networks (ANNs)** are modeled on the human brain
- Nodes have connections ala neurons in human nervous system
- Nodes are also called **neurodes**
- Three types of nodes: input layer, *hidden* layer, output layer



- ANN with one hidden layer is considered as *two-layered*
 - Input layer is not counted
- Can have multiple hidden layers
- Learning through ANNs is also called **connectionist learning**

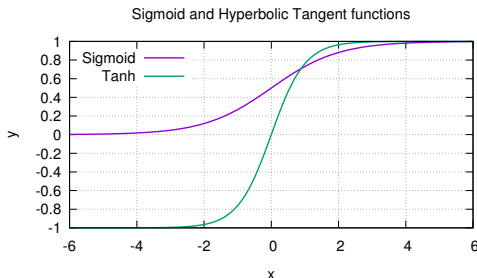
Connections

- ANNs are of many types depending on the connections
- The most common are **multilayer feed-forward networks**
 - Connections are directed from one layer *only* to the next
 - There are *no* back edges or same-layer connections
- In **recurrent networks**, there can be back edges or same-layer connections
- *Fully connected*, i.e., each node in one layer is connected to every node in the next layer
- Training is learning the weights on each of these edges
- Akin to layers of perceptrons
- Activation functions in the nodes are *not* linear
 - Combination of linear functions can only learn linear separators
- Activation function is **sigmoid (logistic)** or **hyperbolic tangent**

Sigmoid and Hyperbolic Tangent Functions

- If input of a node is x , then output y is

$$y = \sigma(x) = \frac{1}{1 + e^{-x}} \quad y = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



- Approximates the step (or sign) function
- Continuous and differentiable
- Output constrained to $(0, 1)$ or $(-1, +1)$
- Scaled versions of each other
- Also called **squashing functions**

Nodes and Weights

- Output of a node is the sigmoid of the weighted sum of its inputs
- Inputs, in turn, are outputs of previous layers
- Inputs are normalized to $(0, 1)$
- Outputs are already constrained to $(0, 1)$
- Weight from a node i to node j is w_{ij}
- Output from node i is O_i
- Input to node j is I_j

$$I_j = \sum_{\forall i} w_{ij} O_i$$

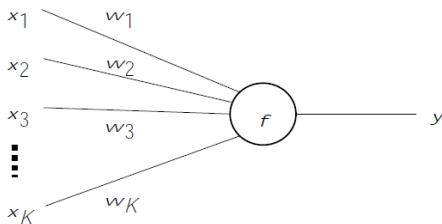
- Output from node j is O_j

$$O_j = \sigma(I_j) = \frac{1}{1 + e^{-I_j}}$$

Training an ANN

- Training an ANN requires
 - Designing the topology: how many layers and many nodes in each layer
 - Learning the weights of the connections
- Designing the topology requires either extensive domain knowledge or simply try-and-test
- The final outputs are some non-linear functions of inputs
- Hence, can be trained using gradient descent
- However, it is too complex and slow
- Weights are updated through the **backpropagation** algorithm

Backpropagation



- Main idea
 - Start with arbitrary weights
 - Propagate forward the values
 - Propagate backward the errors
 - Update the weights using gradient descent

Backpropagation

- Output, using sigmoid function $\sigma(u) = \frac{1}{1+e^{-u}}$, is

$$y = \sigma(u) = \sigma\left(\sum_{\forall x_i} w_i \cdot x_i\right)$$

- Derivatives of activation functions

$$\text{Sigmoid: } \frac{d\sigma(u)}{du} = \sigma(u)\sigma(-u) = \sigma(u)(1 - \sigma(u))$$

$$\text{Tanh: } \frac{d(\tanh(u))}{d(u)} = 1 - \tanh^2(u)$$

$$\text{Linear: } \frac{d(\text{ReLU}(u))}{d(u)} = \begin{cases} 0 & \text{when } u < 0 \\ 1 & \text{when } u \geq 0 \end{cases}$$

Updating Weights

- If true output is t , error is squared error

$$E = \frac{1}{2}(t - y)^2$$

- *Stochastic gradient descent*
- Error function with respect to a single weight w_i

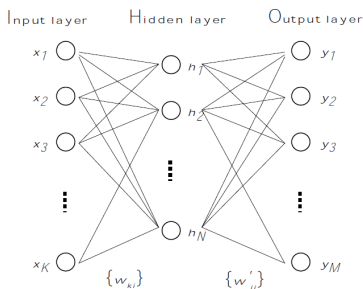
$$\frac{\partial E}{\partial w_i} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial u} \cdot \frac{\partial u}{\partial w_i} = [(y - t)] \cdot [y(1 - y)] \cdot [x_i]$$

- Therefore, update equation is

$$w_i^{(new)} = w_i^{(old)} - \eta \cdot (y - t) \cdot y(1 - y) \cdot x_i$$

- **Learning rate** or *momentum* η controls the speed of training
- Can be continued till error is below a threshold

Backpropagation for Multiple Outputs



$$h_i = \sigma(u_i) = \sigma\left(\sum_{k=1}^K w_{ki} \cdot x_{ki}\right)$$

$$y_j = \sigma(u'_j) = \sigma\left(\sum_{i=1}^N w'_{ij} \cdot x_{ij}\right)$$

$$E(\vec{x}, \vec{t}, W, W') = \frac{1}{2} \sum_{j=1}^M (y_j - t_j)^2$$

Representational Power of an ANN

- Boolean functions
 - Can approximate *any* boolean function with one layer of hidden nodes
 - Number of hidden nodes may be equal to exponential factor of number of boolean variables
 - Output layer wired through AND or OR
- Continuous functions
 - Can approximate *any* continuous function with one layer of hidden nodes up to any arbitrary error factor
 - Due to properties of sigmoid/tanh functions
 - Number of hidden nodes may be large
- Arbitrary functions
 - Can approximate *any* arbitrary function with *two* layers of hidden nodes up to any arbitrary error factor
 - Due to properties of sigmoid/tanh functions
 - Number of hidden nodes may be large

Discussion

- Determining number of nodes and layers is problematic
- **Universal approximators**
 - Sigmoid and hyperbolic tangent functions
- Redundant or irrelevant features will have less weight
- Robust to noise
- Overfitting is a problem
- Gradient descent converges to local minimum only
- Susceptible to class imbalance problem
- Generally, requires large training data
- Very slow to train
- Can be easily parallelized
- Notoriously non-explainable

Term Embedding

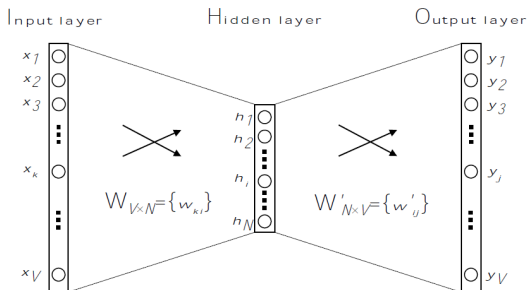
- Aim is to learn an embedding vector for each term that can *predict another term* in its context
- Thus, modeled as a classification problem
- ANNs used: **Word2Vec** model
- If vocabulary size is V , this is a V -class classification problem
- Input is a term *vector*
- Dimensionality is V
- **One-hot vector**
 - Only the specific term is 1, rest are all 0
- Output layer consists of V nodes
- Only one output vector of dimensionality V
- **Softmax** used for classification
- Probability of a particular dimension i with value f_i is

$$p(o_i) = \frac{\exp(f_i)}{\sum_{\forall i} \exp(f_i)}$$

- Choose the word whose probability is the *highest*

One-Word Context

- One word per context is used to predict one target word
 - Subhas : ? Bose



- Only one hidden layer of size N
- Fully connected feed-forward architecture
- Number of weights is $V \times N + N \times V$
- Given a one-hot encoded vector \vec{x}_I for input w_I , the output is \vec{y}_O predicting the word w_O

Layers

- Weight matrix W of size $V \times N$ from input to hidden layer
- Row k represents vector for word k : v_k^T
- For input vector x of word k

$$h = W^T x = W^T I_{(k), \cdot} = v_{w_l}^T$$

- v_{w_l} is the vector representation of input word W_l in N dimensions
- Essentially, the hidden layer is a *linear* copy
- Different weight matrix W' of size $N \times V$ from hidden to output layer
- Score for each word w_j is

$$u_j = v_{w_j}'^T h = v_{w_j}'^T v_{w_l}$$

- v_{w_j}' is the j -th column of W'

Objective

- Using softmax, a multinomial distribution over all predicted words w_j given an input word w_I is defined
- *Log-linear* classification model
- Probability of target word being w_j given an input context word w_I is

$$p(w_j|w_I) = y_j = \frac{\exp(u_j)}{\sum_{j'=1}^V \exp(u_{j'})} = \frac{\exp(v'_{w_j} v_{w_I})}{\sum_{j'=1}^V \exp(v'_{w_{j'}} v_{w_I})}$$

- v_w is the *input vector* of word w
- v'_w is the *output vector* of word w
- Training objective is to *maximize* the probabilities
- Which is the word vector representation?
- The hidden layer of dimensionality N
- Training is through backpropagation
- The output layer is ultimately thrown away!

- Objective

$$\max p(w_{j*}|w_I) = \max y_{j*} = \max \log y_{j*} = u_{j*} - \log \sum_{j'=1}^V \exp(u_{j'}) = -E$$

- Stochastic gradient descent

$$\frac{\partial E}{\partial w'_{ij}} = \frac{\partial E}{\partial u_j} \cdot \frac{\partial u_j}{\partial w'_{ij}} = e_j \cdot h_i$$

$$\frac{\partial E}{\partial w_{ki}} = \frac{\partial E}{\partial h_i} \cdot \frac{\partial h_i}{\partial w_{ki}} = \left(\sum_{j=1}^V \frac{\partial E}{\partial u_j} \cdot \frac{\partial u_j}{\partial h_i} \right) \cdot \frac{\partial h_i}{\partial w_{ki}} = \sum_{j=1}^V e_j \cdot w'_{ij} \cdot x_k$$

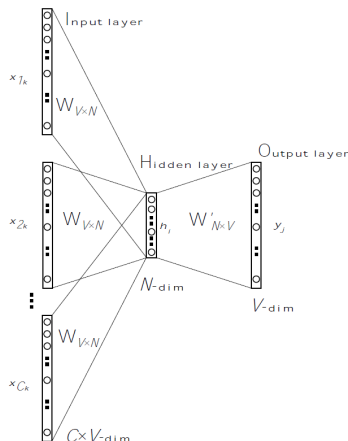
- Update equations

$$\vec{v}'_{w_j}^{(new)} = \vec{v}'_{w_j}^{(old)} - \eta \cdot e_j \cdot \vec{h}$$

$$\vec{v}_{w_I}^{(new)} = \vec{v}_{w_I}^{(old)} - \eta \cdot \frac{\partial E}{\partial h}$$

Continuous Bag-of-Words (CBOW) Model

- *Multiple* context words used to predict a *single* target word
 - Tendulkar, Dravid : ? Ganguly
 - Sachin, Rahul : ? Dev Varman



- Bag-of-words in a local context window that continuously changes

Model

- Instead of just a copy, hidden node is average of C context words

$$\vec{h} = \frac{1}{C} W^T (x_1 + x_2 + \dots + x_C) = \frac{1}{C} (v_{w_1} + v_{w_2} + \dots + v_{w_C})^T$$

- Error

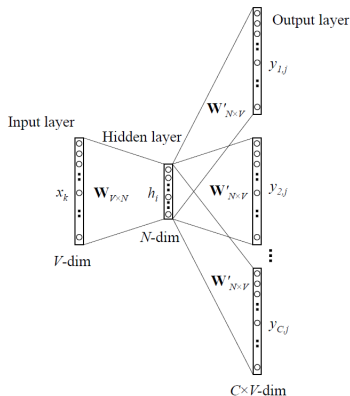
$$E = -\log p(w_{j*} | w_{l,1}, \dots, w_{l,C}) = -u_{j*} + \log \sum_{j'=1}^V \exp(u'_{j'})$$

- Update of hidden layer weights is average

$$\begin{aligned}\vec{v}'_{w_j}^{(new)} &= \vec{v}'_{w_j}^{(old)} - \eta \cdot e_j \cdot \vec{h} \\ \vec{v}_{w_{l,c}}^{(new)} &= \vec{v}_{w_{l,c}}^{(old)} - \eta \cdot \frac{1}{C} \cdot \frac{\partial E}{\partial h}\end{aligned}$$

Skip-Gram Model

- *Single* context word to predict *multiple* target words
 - Dev Varman : ? , ? Sachin, Rahul



- Output layer is C target words of dimensionality V

Model

- Output layer weights are shared, i.e., they are the same matrix W'
- Probability of c^{th} target word being $w_{c,j}$ given an input context word w_l is

$$p(w_{c,j}|w_l) = y_{c,j} = \frac{\exp(u_{c,j})}{\sum_{j'=1}^V \exp(u_{j'})} = \frac{\exp(v'_{w_j}{}^T v_{w_l})}{\sum_{j'=1}^V \exp(v'_{w_{j'}}{}^T v_{w_l})}$$

- Error

$$E = -\log p(w_{j^*,1}, \dots, w_{j^*,c}|w_l) = -\log \prod_{c=1}^C \frac{\exp(u_{c,j^*_c})}{\sum_{j'=1}^V \exp(u_{j'})}$$

- Update of output layer weights is sum

$$\vec{v}'_{w_j}{}^{(new)} = \vec{v}'_{w_j}{}^{(old)} - \eta \cdot \sum_{c=1}^C e_{c,j} \cdot \vec{h}$$
$$\vec{v}_{w_l,c}{}^{(new)} = \vec{v}_{w_l,c}{}^{(old)} - \eta \cdot \frac{1}{C} \cdot \frac{\partial E}{\partial h}$$

Word2Vec Discussion

- Context window size used is 4 words up and down
- Note that, k is chosen randomly within 4
 - This is why it is called “continuous” bag-of-words
- Number of context words, C , is within 5
- Rare words are discarded
 - This also effectively increases context window size
- Log-linear model
 - This makes training faster
- Input word and context word are treated differently
 - No good reason except mathematical convenience
- Skip-gram gives better semantic results (poorer syntactic, and marginally better overall)
- Corpus is more important than method
- Why does just context perform well?
 - No definitive or satisfying answer

Training

- Update equations require V negative training equations per word
- This is prohibitively expensive
- **Hierarchical softmax** models a binary tree with the V words as leaves
- Training or reaching a word vector requires $\log_2 V$ operations
- **Negative sampling** updates only a sample of words instead of all V
- Positive word must be there in the sample
- Some more “negative” words (around 10 suffices!)

$$P(w_j|w_t) = \log \sigma(\tilde{v}_j^T v_t) + \sum_{i=1}^k \mathbb{E}_{w_i \sim P_n} \log \sigma(-\tilde{v}_i^T v_t)$$

- Sampled using probability distribution $p(w) \sim P_n \propto f_w^{3/4}$
- **Sub-sampling** of frequent words
 - Essentially, stopwords; e.g., “the” is in the context of almost every noun
- Word w is discarded with probability $p(w) = 1 - \sqrt{\frac{t}{f_w}}$ using a small threshold $t \sim 10^{-5}$

Morphologically Richer Languages

- Morphologically richer languages have the same internal root or structure for a large number of words
- Further, it requires a prohibitively large corpus to train for all such variants of a word
- Examples: Indian languages, Finnish, Turkish, and even European languages such as French, Spanish, German, Russian
 - In Samskritam, from a single *dhatu*, around 1 lakh forms can be produced (not all of them are words, though)
 - With around 2,000 dhatus, number of forms is 200 million
- Simple corpus-based word vector embeddings may not work
- **FastText** uses **character n-grams**
- Originally called **SISG (Subword Information Skip Gram)**

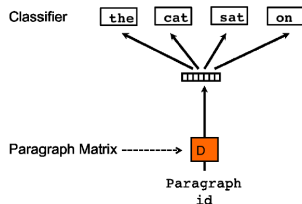
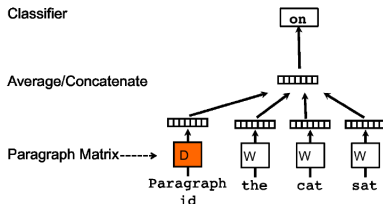
- A word vector representation is sum of character n-gram vectors and the word vector itself
- Two limits for n-gram length: minimum and maximum
 - Generally, 3 and 6
- Example word: “India”
- If limits are 2 and 3, then embeddings are of “In”, “Ind”, “ndi”, “dia”, “ia”, and “India”
- If G_w is the set of n-grams for a word w , then

$$s(w, c) = \sum_{\forall g \in G_w} z_g^T v_c$$

- **FastText** can handle *out-of-vocabulary* words
 - Glove and Word2Vec cannot!
- Uses as many common n-grams as it can from the new word

Document Vectors

- Vectors for sentences, paragraphs, and documents
- Average of word vectors
- Can be directly learnt using an architecture similar to Word2Vec



- Paragraph vectors (also) contribute to predicting the context word
- Paragraph vector is shared *only* within a paragraph
- Paragraph vectors are of fixed dimensionality K
- Model 1: **Distributed Memory Model of Paragraph Vectors (PV-DM)**
 - Word vectors are shared across the corpus
 - Total number of parameters is $P \times K + V \times N$
- Model 2: **Distributed Bag-of-Words version of Paragraph Vectors (PV-DBOW)**
 - Total number of parameters is $P \times K$