

Text Classification and Neural Networks

Jelke Bloem & Giovanni Colavizza

Text Mining
Amsterdam University College

March 21, 2022

Announcements

- Assignment 2 due tomorrow
- Next reading assignment in two weeks
- After Assignment 2, think about project topics and groups

Overview

- 1 Text Classification
- 2 Logistic Regression
- 3 Evaluation
- 4 Neural Networks
- 5 Extras

Word Embeddings reading assignment: Questions

- How is Shifted PPMI actually shifted and what advantage does that give it over standard PPMI?
 $PMI(w, c) \log(k)$ where k = number of negative samples, constant
- Why are we using models that we do not totally understand how they work and why they produce effective results?
- How does the factorization of a matrix affect our word vectors?
- How can some models be simpler and be better than more complex models?
- What are Large Scale Distributed Deep Networks such as DistBelief, what makes them different from standard neural networks (other than their size)?
- How exactly is word similarity decided? For example in a way “big” and “small” are similar in that they both talk about size, but at the same time are complete opposites when it comes to the scale of something.
- How does the Huffman binary tree store data more efficiently?

Word Embeddings reading assignment: Questions

- I wonder how different evaluation metrics influence the results. Is an “improved” model only improved under the parameters of a certain evaluation method? How could this be generalized?
- Mikolov et al. suggested in their paper that their model is more time-efficient and can train with larger datasets than neural network approaches due to the lower computational complexity. I wonder if it is still true?
- Can we utilise CBOW and Skip-gram on other tasks such as speech recognition and automated essay scoring? Can SVD be used in different fields, such as recommendation systems?
- What does the implicit matrix factorization indicate in real-life?
- Since this paper is from 2013 I was wondering how much development happened in the time until now and if those models have been replaced with better ones or how they have evolved.
 - ▶ Transformer models / contextual word embeddings
- For the linguistic tasks in the data set values for word similarities were assigned by humans. Does that not pose a problem since people can just estimate similarity but there is not really a value you can assign to it?

Word Embeddings reading assignment: Questions

- If computation efficiency increases (from a hardware point of view), would at some point the models including a non-linear hidden layer become better at this task again?
- If two models are trained on different corpuses but used in the same task how comparable are their performances?
- How does the “noise-contrastive estimation” (NCE) technique differ from SGNS?
 - ▶ Avoids normalizing probabilities over the entire vocabulary. Logistic regression classifier discriminating between data distribution and ‘noise’ distribution
- What is D in Mikolov et al.?
- Is it always beneficial and necessary to use larger datasets in research related to NLP?
- What constitutes accuracy in a word similarity task?
 - ▶ Mikolov et al.: Question answering task
- In the skip-gram method, is it possible to treat the entire given corpus as input and calculate the similarities based on that corpus instead of sentence by sentence?

Text Classification

Task definition

- We are given a **training set** $\{X, Y\}$ of data pairs (x, y) , where x is a text document and y is the class the document belongs to.
- Each $y \in \mathcal{Y}$, where $\mathcal{Y} = \{c_1, c_2, \dots, c_k\}$ are the distinct (finite and enumerable) classes we have. If $|\mathcal{Y}| = k = 2$, we have a binary classification task.
- Using a *learning method*, our goal is to learn a **classifier**, or a classification function γ that maps documents to classes:

$$\gamma : \mathcal{X} \rightarrow \mathcal{Y}$$

- The fact that we use annotated data to learn makes this a form of *supervised learning*. Note that a “document” can be anything really: words, text sequences, longer texts.

Examples

task	x	y
language ID	text	{english, mandarin, greek, ...}
spam classification	email	{spam, not spam}
authorship attribution	text	{jk rowling, james joyce, ...}
genre classification	novel	{detective, romance, gothic, ...}
sentiment analysis	text	{postive, negative, neutral, mixed}

Credit: David Bamman (UC Berkeley).

Text representation

Our text documents X can be **represented** in many ways:

- Pre-computed features (e.g., the length of the document or the average length of the words it contains).
- A selection of words (e.g., only stopwords for language detection).
- Words in isolation (so called “bag of words”, or unigram model).
- Conjunctions of words (e.g., bigrams).
- Higher-order features (e.g., PoS).
- Word embeddings.

Logistic Regression

Logistic regression

- Our goal is, given a document represented with a feature vector \mathbf{x} and classes $c \in \mathcal{Y}$, to learn a classifier discriminating the right class for \mathbf{x} :

$$\hat{p}(y = c|\mathbf{x})$$

- Let us start with a binary classifier and two classes, thus $\mathcal{Y} = \{0, 1\}$.
- We need to estimate $\hat{p}(y = 1|\mathbf{x})$, and $\hat{p}(y = 0|\mathbf{x}) = 1 - \hat{p}(y = 1|\mathbf{x})$ will follow suit.
- Logistic regression uses two components for this: a **linear model** of the inputs and the **Sigmoid (or logistic) function**. So, it is like the perceptron but with a different classification function.

Sigmoid (or logistic) function

- Let us consider the set of features x_1, x_2, \dots, x_d we used to represent our input document \mathbf{x} . We add $x_0 = 1$ to model the intercept, and create a linear model with them:

$$z = \sum_{j=0}^d w_j x_j = \mathbf{w} \cdot \mathbf{x}$$

- To create a probability distribution, we pass z through the Sigmoid $\sigma(z)$:

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

- The Sigmoid squeezes z within 0 and 1 and is always positive.

Sigmoid (or logistic) function

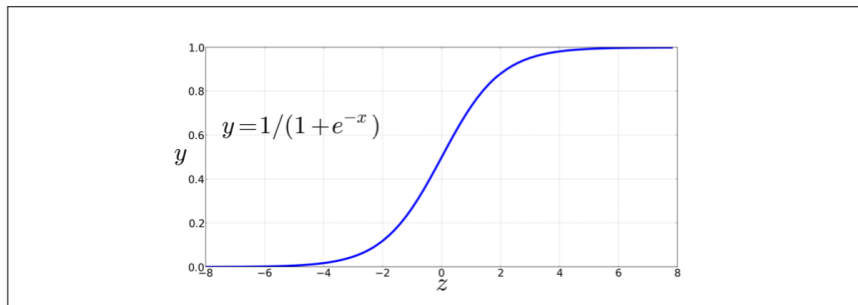


Figure 5.1 The sigmoid function $y = \frac{1}{1+e^{-x}}$ takes a real value and maps it to the range $[0, 1]$. Because it is nearly linear around 0 but has a sharp slope toward the ends, it tends to squash outlier values toward 0 or 1.

Credit: M&J, Ch. 5.

Logistic regression

- Applied to our binary classification task, we have that:

$$\hat{p}(y = 1|\mathbf{x}) = \sigma(z_{\mathbf{x}}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$
$$\hat{p}(y = 0|\mathbf{x}) = 1 - \sigma(z_{\mathbf{x}}) = \frac{e^{-\mathbf{w} \cdot \mathbf{x}}}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

- Then, we just need to use a **decision boundary** to assign the class given the estimated probabilities:

$$\hat{y} = \begin{cases} 1 & \text{if } \hat{p}(y = 1|\mathbf{x}) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

- So, we have defined our data and task, and have a model.
What do we miss?

Logistic regression: Cross-entropy

- We need a loss function. Let us use MLE to find one.
 - ▶ Maximize the conditional probability of observing the data given a distribution
- We have that $p(y|\mathbf{x})$ follows a Bernoulli distribution given that we only have two discrete outcomes $(0, 1)$, hence:

$$p(y|\mathbf{x}) = \hat{y}^y(1 - \hat{y})^{1-y}$$

- As usual, let us move to log space and add a minus to switch to a minimization problem (note we work with a single data point (\mathbf{x}, y) for now):

$$\begin{aligned} -\log p(y|\mathbf{x}) &= -\log [\hat{y}^y(1 - \hat{y})^{1-y}] \\ &= -[y\log \hat{y} + (1 - y)\log(1 - \hat{y})] \end{aligned}$$

- Let us now plug-in the Sigmoid and call it the loss:

$$\mathcal{L}_{\mathbf{x}}(\mathbf{w}) = -[y\log \sigma(\mathbf{w}\mathbf{x}) + (1 - y)\log(1 - \sigma(\mathbf{w}\mathbf{x}))]$$

Logistic regression: Cross-entropy

- Let us now plug-in the Sigmoid and call it the loss:

$$\mathcal{L}_x(\mathbf{w}) = -[y \log \sigma(\mathbf{w}\mathbf{x}) + (1 - y) \log(1 - \sigma(\mathbf{w}\mathbf{x}))]$$

- The loss on the whole dataset is going to be (note we are already in log space thus we can sum):

$$\mathcal{L}(\mathbf{w}) = -\frac{1}{N} \sum_{i=1}^N [y_i \log \sigma(\mathbf{w}\mathbf{x}_i) + (1 - y_i) \log(1 - \sigma(\mathbf{w}\mathbf{x}_i))]$$

- Equivalent to calculating cross-entropy for the Bernoulli distribution
- To this we can, as usual, attach regularization:

$$\mathcal{L}_{L_2}(\mathbf{w}) = \mathcal{L}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Logistic regression: Optimization via SGD

- The last missing bit is how to find good parameters \mathbf{w} : we can use SGD.
 - ▶ There is no analytical solution, unlike linear regression
- It turns out that the derivative for one data point \mathbf{x} is (w.o. regularization):

$$\frac{\partial \mathcal{L}_{\mathbf{x}}(\mathbf{w})}{\partial \mathbf{w}_j} = [\sigma(\mathbf{w}\mathbf{x}) - y] \mathbf{x}_j$$

- For multiple data points, we just sum (w.o. regularization), and with this we are good to go for SGD:

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}_j} = \sum_{i=1}^N [\sigma(\mathbf{w}\mathbf{x}_i) - y_i] \mathbf{x}_{ij}$$

- *Full derivation as an extra, below.*

Evaluation

Data splitting

	training	development	testing
size	80%	10%	10%
purpose	training models	model selection; hyperparameter tuning	evaluation; never look at it until the very end

Credit: David Bamman (UC Berkeley).

Accuracy and baselines

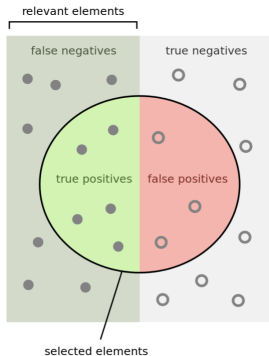
- **Accuracy** is the fraction of correctly predicted data points over the total. It can be calculated on any dataset split: train, development and test. Very good starting point.
- **Baseline**: important to have one. It can be a random classifier (i.e., flip a coin for a binary classifier), or a fast and reasonable model (e.g., logistic regression with TF-IDF features).

Precision and recall

Given a binary classifier:

- **True positive:** a data point correctly predicted to be 1.
- **True negative:** a data point correctly predicted to be 0.
- **False positive:** a data point incorrectly predicted to be 1.
- **False negative:** a data point incorrectly predicted to be 0.

Precision and recall



How many selected
items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant
items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Credit: Wikipedia.

F-measure and accuracy reloaded

- F-measure (harmonic mean of precision and recall):

$$F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

- Accuracy:

$$A = \frac{tp + tn}{tp + tn + fp + fn}$$

Parameters and hyperparameters

Parameters whose values are *learned*

Feature	β
the	0.01
and	0.03
bravest	1.4
love	3.1
loved	1.2
genius	0.5
<i>BIAS</i>	-0.1

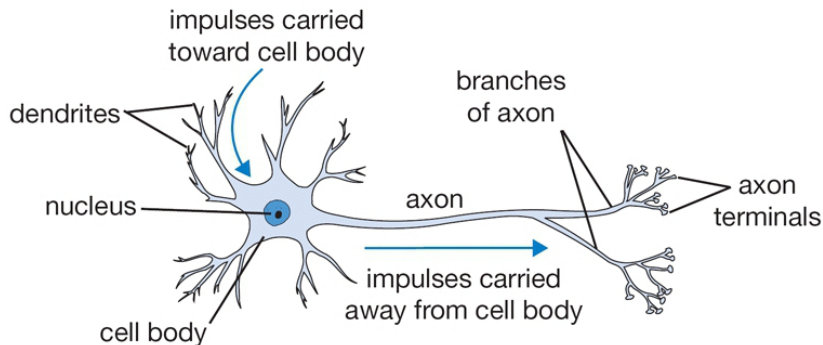
Hyperparameters whose values are *chosen*

Hyperparameter	value
minimum word frequency	5
max vocab size	10000
lowercase	TRUE
regularization strength	1.0

Credit: David Bamman (UC Berkeley).

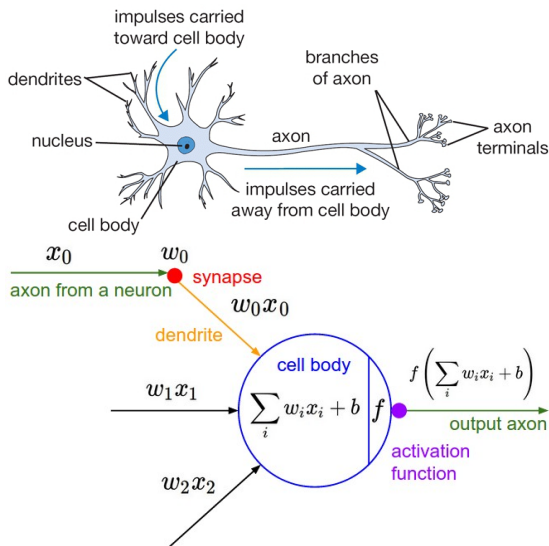
Neural Networks

A single neuron



Credit: Andrej Karpathy via Stanford's CS231N.

A single neuron



Credit: Andrej Karpathy via Stanford's CS231N.

Logistic regression as a neural network

Following the notation in the previous slide, we have:

- $\mathbf{x} = \langle x_0, x_1, x_2, \dots, x_d \rangle$ is our input representation.
- We aggregate the features \mathbf{x} into a linear combination using weights \mathbf{w} . We also include the bias term b into the matrix by adding an appropriate dimension fixed at 1 to \mathbf{x} , so that we can use matrix notation:

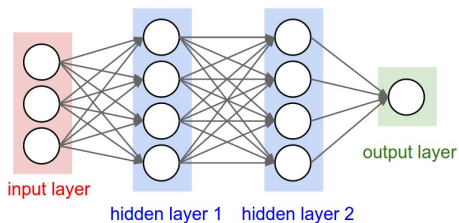
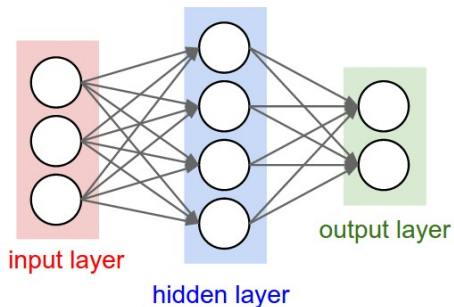
$$z = \sum_{i=0}^d w_i x_i = \mathbf{w} \cdot \mathbf{x}$$

- We pass z through an activation function, in this case the sigmoid:

$$f = \sigma(z) = \frac{1}{1 + e^{-z}}$$

- **Linear models are a single neuron.** *Question: what is the activation function for linear regression?*

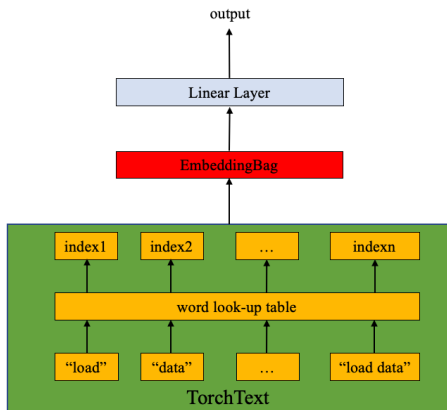
From single layer to multi-layer



Credit: Andrej Karpathy via Stanford's CS231N.

Using embeddings as features

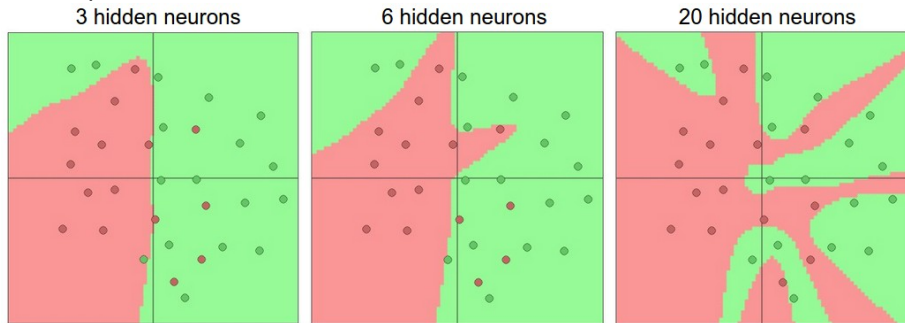
Neural networks are **modular**: we can piece them together into advanced architectures. For example, we can use embeddings to represent our input (either training them or using pre-trained ones). *More on this in the lab.*



Credit: TorchText.

Why do we need non-linearities?

Multiple layers and **non-linear functions** (such as the sigmoid) allow us to fit complex decision boundaries.



Credit: Andrej Karpathy via Stanford's CS231N.

How do we train neural networks?

- Key idea: use a smart way to apply SGD, called **backpropagation**.
- Backpropagation combines using the chain rule to calculate local derivatives (called gradients) with the re-use of pre-computed operations to speed the computation up.
- *More on this in the external materials for the course.*

Neural networks practicalities

Training neural networks entails a lot more than stacking up layers. Several topics require practical and theoretical knowledge beyond this course:

- Weight initialization
- Regularization (e.g., via dropout)
- Which non-linearities to use
- Which loss functions to use
- How to monitor and adjust the learning process (e.g., optimizers and learning rates) to avoid dying neurons and overfitting

Extras

Full derivation for logistic regression

- First, we need some notable derivatives:

$$\frac{\partial \log(x)}{\partial x} = \frac{1}{x}$$

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$$

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x} \rightarrow \text{chain rule}$$

Full derivation for logistic regression

- Then:

$$\begin{aligned}\frac{\partial \mathcal{L}_x(\mathbf{w})}{\partial w_j} &= -\partial [y \log \sigma(\mathbf{w}\mathbf{x}) + (1 - y) \log(1 - \sigma(\mathbf{w}\mathbf{x}))] \\ &= -[\partial y \log \sigma(\mathbf{w}\mathbf{x}) + \partial(1 - y) \log(1 - \sigma(\mathbf{w}\mathbf{x}))] \\ &= -\frac{y}{\sigma(\mathbf{w}\mathbf{x})} \partial \sigma(\mathbf{w}\mathbf{x}) - \frac{1 - y}{1 - \sigma(\mathbf{w}\mathbf{x})} \partial(1 - \sigma(\mathbf{w}\mathbf{x})) \rightarrow \text{chain rule} \\ &= -\left[\frac{y}{\sigma(\mathbf{w}\mathbf{x})} - \frac{1 - y}{1 - \sigma(\mathbf{w}\mathbf{x})} \right] \partial \sigma(\mathbf{w}\mathbf{x}) \rightarrow \text{re-arrange}\end{aligned}$$

- *Exercise: plug-in the derivative of the Sigmoid and re-arrange yourself to reach:*

$$\dots = [\sigma(\mathbf{w}\mathbf{x} - y)] x_j$$

Full derivation for logistic regression

- In case you were wondering:

$$\begin{aligned}\frac{\partial \sigma(x)}{\partial x} &= \partial \frac{1}{1 + e^{-x}} \\ &= \partial [1 + e^{-x}]^{-1} \\ &= \frac{e^{-x}}{1 + e^{-x}} \frac{1}{1 + e^{-x}} \\ &= \frac{(1 + e^{-x}) - 1}{1 + e^{-x}} \sigma(x) \\ &= \sigma(x)(1 - \sigma(x))\end{aligned}$$

- *Exercise, derive:*

$$\frac{\partial \log \sigma(x)}{\partial x} = \sigma(-x)$$

Why MSE and cross-entropy?

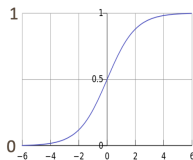
- It turns out that, given some standard assumptions on our models, using those two losses corresponds to doing Maximum Likelihood Estimation. See <https://www.expunctis.com/2019/01/27/Loss-functions.html>.
- If you are curious about the information theory underpinning cross-entropy, read this: <http://colah.github.io/posts/2015-09-Visual-Information>.

NN activation functions

Several non-linear activation functions have been proposed. A good default options is ReLU.

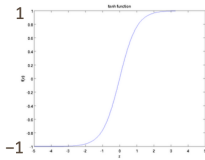
logistic ("sigmoid")

$$f(z) = \frac{1}{1 + \exp(-z)}.$$



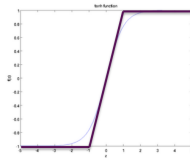
tanh

$$f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}},$$



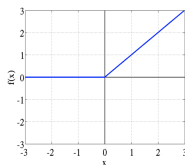
hard tanh

$$\text{HardTanh}(x) = \begin{cases} -1 & \text{if } x < -1 \\ x & \text{if } -1 \leq x \leq 1 \\ 1 & \text{if } x > 1 \end{cases}$$



ReLU (Rectified Linear Unit)

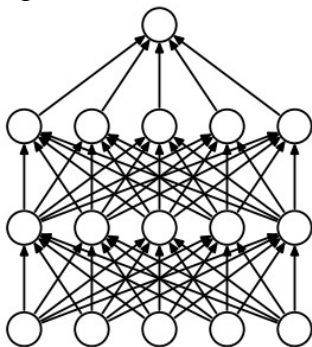
$$\text{rect}(z) = \max(z, 0)$$



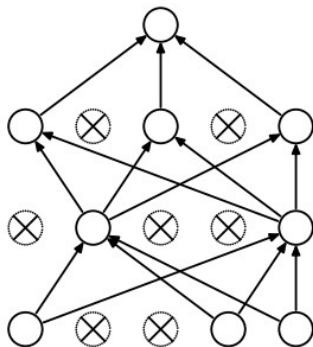
Credit: Stanford CS224N.

NN regularization via dropout

Dropout's idea is to mask a random set of neuron connections at training time, in order to compel the network to learn redundant paths and avoid overfitting.



(a) Standard Neural Net



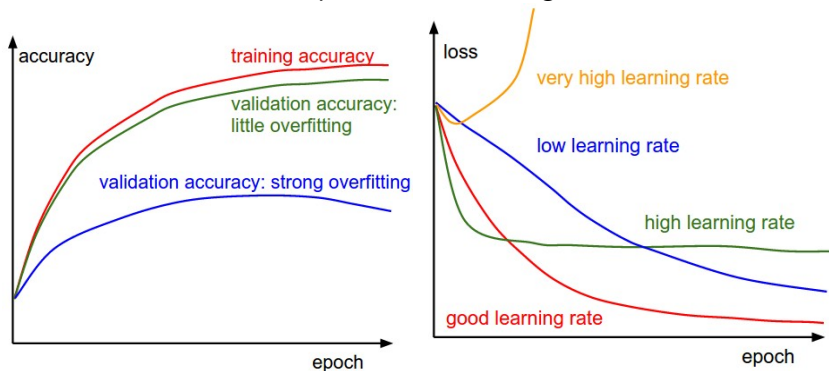
(b) After applying dropout.

Credit: Srivastava et al.

<https://www.cs.toronto.edu/~hinton/absps/JMLRdropout.pdf>.

NN under/overfitting and learning rates

Two illustrations on how to spot correct learning behaviour.



Credit: Andrej Karpathy via Stanford's CS231N.