

ECE 1786 Lecture #3

Work-in-Flight: Assignment 2 - Classification of Language, due Monday Oct 7

Next Week: Will discuss the structure and timing of the course project; do not form teams or discuss in any detail the project yet; **hyperparameter N not set**

Last Day: How Word Embeddings are Trained

Today: Classification of Language using word embeddings

Assignment 1 due last night, latest to hand in is tonight at 9pm, with penalty

How to keep up with this field? Two good Media to keep your eyes (& ears) on:

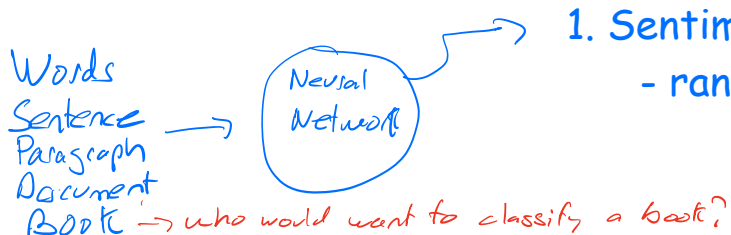
1. "The Batch" a weekly newsletter from Andrew Ng, and AI luminary:
 - <https://www.deeplearning.ai/the-batch/>
 - very interesting, timely, tightly written about AI, much about LLMs
2. The Practical AI Podcast - <https://changelog.com/practicalai>
 - Two practitioners talking about/interviewing on many interesting topics
 - originally ML/Data science; now lots on LLMs, generative AI

After Assignment 1, we now know/see how embeddings respect meaning

- Two words with the same meaning would have similar embeddings, so you don't have to have a program deal with specific words
- But: One word with multiple meanings: still a problem -> later

A key general application in machine learning is classification

- You have seen picture classification in your prior ML class(es)
- Now, we'd like to classify words, sentences, paragraphs and documents
 - Into what, though?



2. Named Entity Recognition: identify thing that can described in many ways
 - e.g. recognize a specific reason for quitting smoking: Name: Calming = "relaxes me" or "makes me calm" or "feeling smooth"
 - e.g. recognize a noun or verb or adjective (quite general)
 - e.g. Proper name
3. Psychology: identify a specific style of talking from therapeutic perspective [change v. sustain]
4. Politics - left vs. right wing style speech
5. Mental Health: detect Depression/Anxiety in Speech.
6. Law: looking for specific facts for specific facts in a document
 - looking for "privileged" information in emails (information protected)
- Now that we have text represented as embeddings, we can use a Neural Network to work with it/classify it

In Assignment 2, you'll be training two types of networks to detect if a sentence is either 1. Objective (a statement of fact) or 2. Subjective (an opinion).

- You'll also look inside the network to see what it is learning!
 - Dataset used for training comes from Pang & Lee @Cornell who created the two parts of the dataset:
1. The **subjective** sentences come from movie reviews, assumed to be subjective (but may not be)
 - EXAMPLE: "His performance lacked energy"

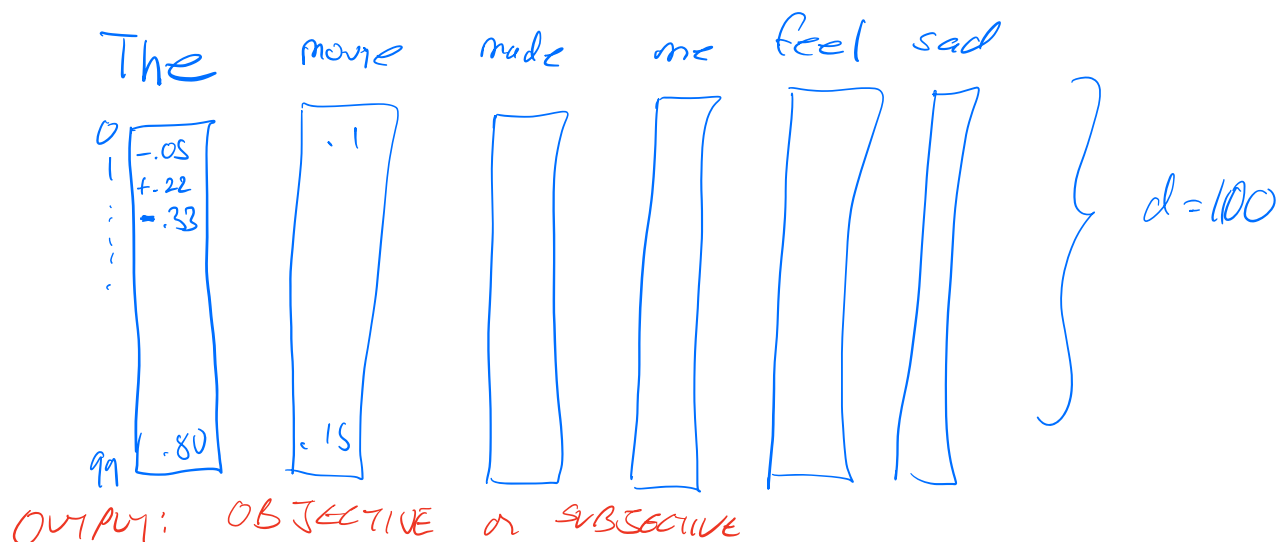
2. The **objective** sentences are taken from plot summaries from IMDB (the Internet Movie Database) - assumed to be objective, but also may not be

- EXAMPLE: "The heroine jumped from the moving train"

- What about: "The movie made me feel sad"? Objective or Subjective?

So, we want to classify sentences as objective or subjective

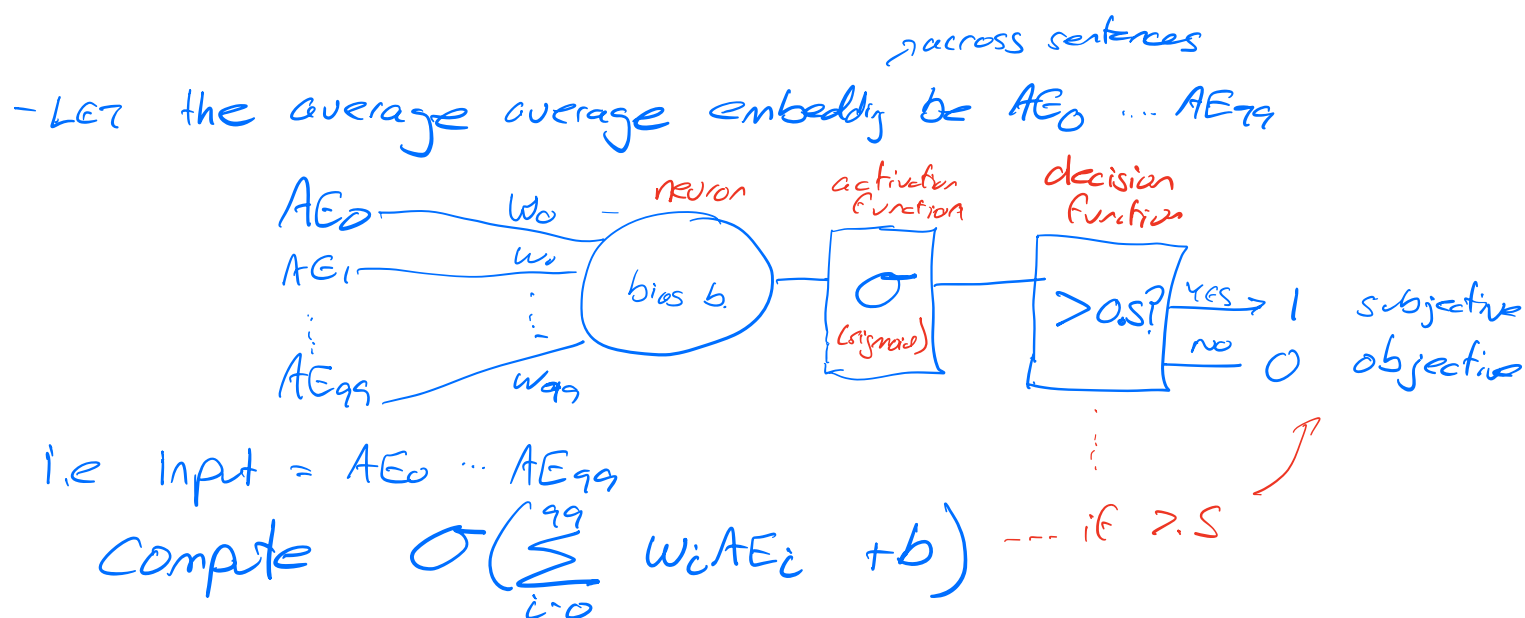
- each of these sentences is input as regular (ASCII) text;
- the first processing step is called 'tokenization' which often (but not always) breaks words into smaller **units**.
 - The lemmatization in you did in A1 was a little like this
- Each of those **units** will need an embedding;
 - in the case that there is a word that doesn't break down into known units, it will be assigned the 'unknown' token and associated embedding
- Assignment 2 uses the GloVe embeddings, as you used in assignment 1 sections 1 and 2, but with dimension **d = 100** this time
- So the input to the models that you'll build in assignment 2 is a sentence of words that is converted into a sequence of embeddings:



- Notice that input sentences would be of different lengths; but some neural nets are set up for fixed-length inputs
 - In that case, must "pad" inputs with zeroes (i.e. embeddings that are zero) to make a batch of all equal size length inputs. Why zeroes?
 - This applies to the CNN, (Method 2 below) but not Method 1

METHOD 1: Single Neuron

- the first model you'll create and train will be a single neuron
 - We'll use this as a baseline to compare to method 2
- The computation of the neural network will be:
 - Compute the average embedding across all words - **what does this accomplish? What could it mean to "average" meaning?**
 - Gives a fixed-length vector/embedding of size 100
 - Feed the result into a single neuron:



- You'll train this neural network - the single neuron - with the given dataset
- It works surprisingly well!
- In other applications that process sentences, the average was commonly used
- **Very interesting:** notice that the trained weights - $w_0 \rightarrow w_{99}$ is a vector of the same size as the embedding. Does it have meaning in the same 'space' as the embedding space? You're asked to explore this in Assignment 2. **How?**

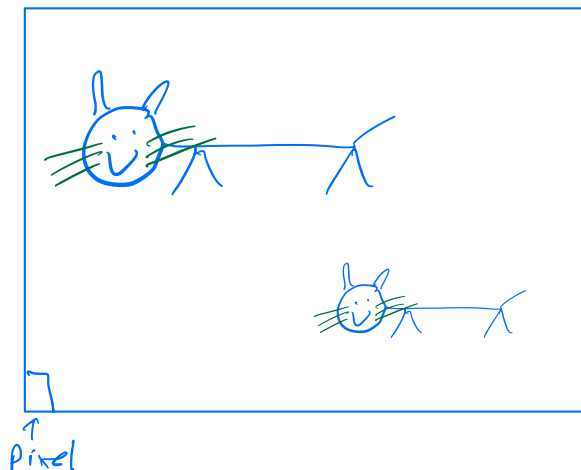
METHOD 2: Convolutional Neural Net (CNN)

Let's first spend a few minutes reviewing CNNs

Recall that pictures are made up of pixels, each pixel is possibly three numbers, the amount of Red, Green and Blue in the pixel

- Say a picture is 1000x1000 pixels:

picture .



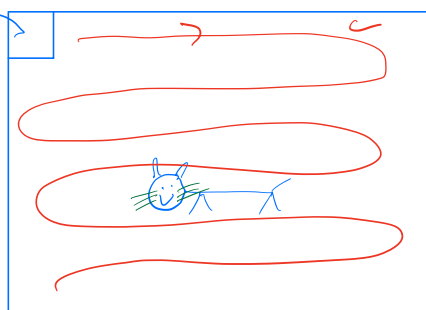
CAT.(S)

Kernel

looks for a specific pattern
e.g.

0	0	1
0	1	0
1	0	0

- looks for diagonal line

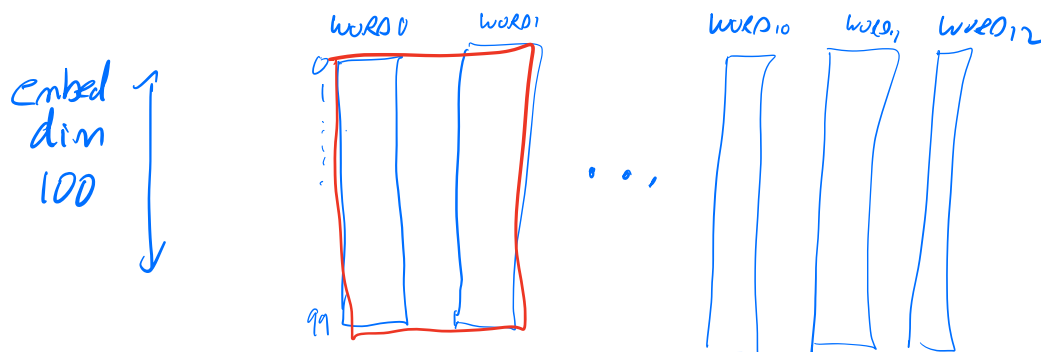


- Kernel scans across the picture "looking" for the smaller picture that is in the kernel.

- the output from a kernel is a new array/picture feature map
- it is the dot product (convolution) of the kernel with each 3x3 sub-image of the image.
- dot product \equiv convolution is a way to detect if the kernel "matches" the sub-image.
- higher dot product means greater match.
 a single value.

- during training, the kernels are learned
- there are lots of kernels.

Now, back to the language classification problem; consider the input sentence(s) embedding sequence:



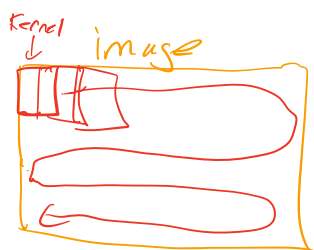
- We would like to look for

- 1) Single words that indicate subjective/objective

His performance lacked energy;

- 2) Pairs of words, triplets, 4 words?

- In general to look for K words
- In CNN-terminology, we'd say that we want to train Kernels of size $K \times 100$
- In a similar way that CNN kernels "scan" (sweep) across the field of a picture (in computer vision), these a kernel would 'scan' across a sentence



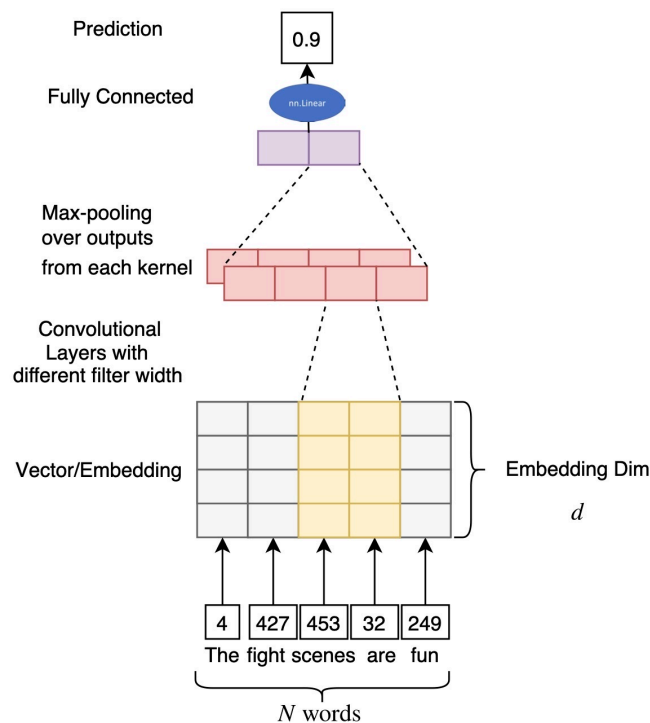
→ kernel is multiplied + accumulated (convolved) with image - dot product
 → remember that during training, the kernel is learned

- In the context of classifying the language sentence, a kernel, which might be of size 2×100 , would just scan across the sentence's word embeddings once.
- It would be trained to look for a 2-word pattern of meaning that would contribute to learning the labels: subjective or objective
- Questions:
 - Is one kernel enough? Maybe not.
 - What size(s) should K be?

- Assignment 2 suggests having n_1 kernels of size $k_1 \times 100$ and n_2 of size $k_2 \times 100$
 - Each one randomly initialized, as all weights/kernels are before training

If you have an input sentence of N words a kernel size of $k_1 = 2 \times 100$? (Where 100 is the embedding dimension?); assume the stride=1. (what is stride?) **What is the size of the output?**

- Get $N-1$ values out; if you have n_1 such kernels, then you get $n_1 \times (N-1)$
- In general for kernel of size k , you get $N-k+1$ values



Yoon & Kim (paper referenced in A2) suggest choosing the maximum across all $N-k+1$ values

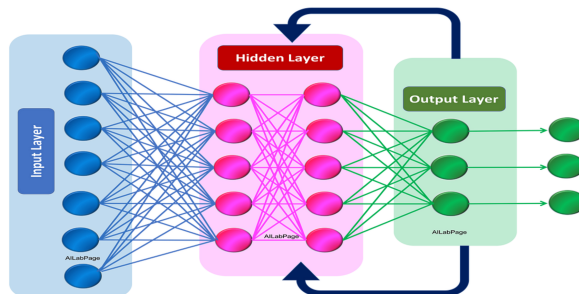
- i.e. maxpool
- Then feed all of those maximum values, from all the kernels into a multi-layer perceptron (MLP) - also called linear, fully connected layer(s).

Does this make sense?

- Should look for patterns of 1-6 words that give sense of whether sentence is objective or subjective

- How might you know what the kernels are trained to look for? i.e. what do they learn? (Discuss)
 - You're asked to explore this in Assignment 2, Section 5.3
 - These should be textual features just as a vision CNN has visual features in its kernels
- **Important note in A2:** for Method 2, I also asked you to enable the training of the embeddings themselves. That is, the gradient descent is set to propagate back into the embeddings, so that they learn to do this task specifically well
 - This is the case for the next network - Transformers - that we'll cover
 - Looking ahead: it will be important to understand that the embeddings themselves are essentially network parameters, and so it makes sense to train them all together
 - the only difference is that you use different parameters depending on what the actual input is

-
- **Aside:** Recurrent Neural Networks (RNNs) had traditionally been used for Natural Language Processing
 - An RNN typically takes 1 embedding in at a time, and has a cycle in which the output feeds back into the network, along with the next input:



- RNNs (LSTMs, GRUs) were always rather problematic in that convergence of training was difficult to achieve and unreliable.

- Also, to me, the fact that all information was “pinched” through the size of the embedding meant that lots of information was lost
- The next (and central) topic of this course is the **Transformer** neural network, which keeps much more information 'alive' in parallel
- Transformers **are closer** to CNN's in that sense, and they have essentially replaced RNNs for NLP

Note on Project: don't use any old RNNs like LSTMs or GRUs; they really have been surpassed. (Unless you've got a really good reason that you check with me)