

ECE 1786 Lecture #3

Last Day: How Word Embeddings are Trained.

Work-in-flight: Assignment 2 Classification

Today: Classification using word embeddings

- Comments on Assignment 1?

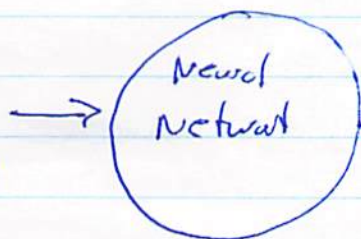
So, now we know/see how embeddings represent meaning

- two words with the same meaning would have similar embeddings
- one word with multiple meanings: still a problem \rightarrow later

A key general application in machine learning is classification.
^{you} have previously seen picture classification

- now we'd like to classify sentences, paragraphs, documents
- into what?

SENTENCE
PARAGRAPHS
DOCUMENTS
BOOK



① Sentiment - positive or negative
 \rightarrow range $-1 \rightarrow +1$

② Named Entity Recognition
 - identify something that can describe in many ways

④ Politics - left v. right.

- e.g. specific reason for getting something
 "makes me calm" \equiv "relaxes me"

⑤ Depression / Anxiety
 in speech.

③ Style of talking
 "charge talk" } in behavior
 "sustain talk" } change

⑥ Suicidality \rightarrow Facebook detects. (Google to see)

⑦ Lawyers \rightarrow look for specific facts/terms

- now that we have text \rightarrow embeddings, we can use a NN to work with it/classify it.

In assignment 2, you be training two types of networks to detect if a sentence is objective (a statement of fact) or subjective (an opinion).

DEMO w. gradto

- you'll also look inside the networks to see what it is learning, (and ~~how~~ use a software tool to gain insight on particular predictions.)

- Dataset used for training: Pang & Lee (Cornell) ^{criteria the 'subjective'}
 SUBJECTIVE SENTENCES: movie reviews from Rotten Tomatoes
 - assumed to be subjective (not always true)

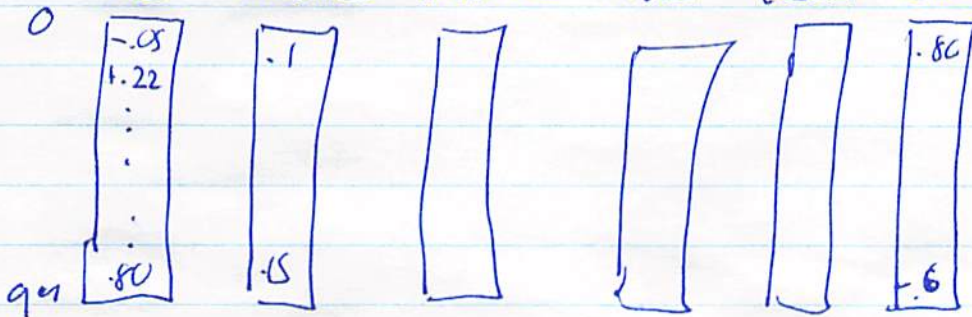
OBJECTIVE SENTENCES: plot summaries from ~~IMDb~~ ~~TMDB~~
 - assumed to be objective (statements of fact about what happened in movies) (not always true)

- do no special tokenization: ^{just spaces} always have "unknown" token for words not in embedding matrix

- will use GloVe embeddings. again, but larger dimension dim=100.

- so the input to the models in A2 is ~~the~~ a sequence of word embeddings.

The movie made me feel sad.

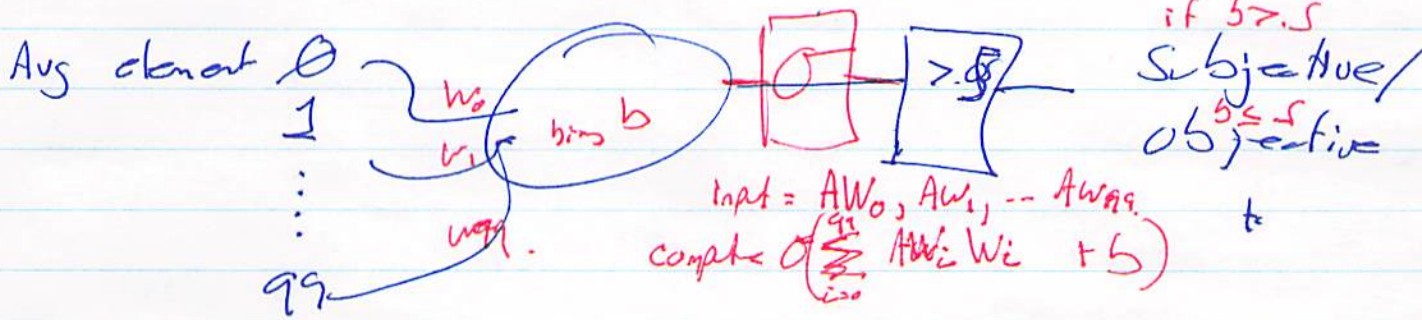


a sentence of 6 words that is converted into

- sentences are of different lengths, while neural nets (except RNNs) are mostly set up for fixed-size inputs.
 → ~~can~~ "pad" inputs with zeroes if ~~necessary~~ ~~inputs~~.
 (~~will be for error~~) to make a batch all equal length

→ discuss.

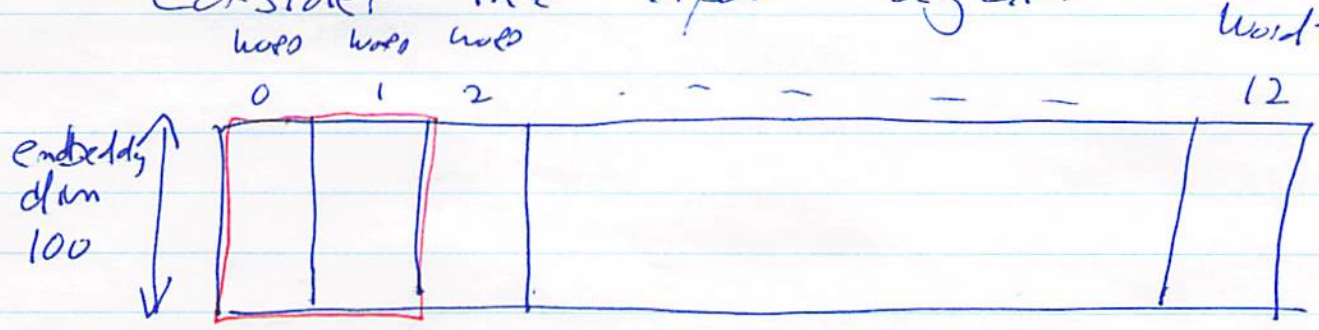
- A baseline model is ~~used~~ first is: 12 0
 - compute the average across all word embeddings - a little in ≈ 1 - recall $\frac{one + ten}{2} \approx five$
 - what does this accomplish → "compressing meaning"
 - get a fixed-length vector. $\frac{happy + sad}{2}$
 - perhaps 100 dimensions is enough to encode ^{meaning}
 - feed result into an MLP (multi-layer perceptron) of just 1 neuron (just 1 layer)



- works surprisingly well
- average is often used to represent sentences.
- $w_0 - w_{99}$ same size as embedding → can explore what it means. Review CNNs

Method 2: A Convolutional Neural Net (CNN)

- Consider the input again



- would like to look for ① Single words that initiate subj. obj. 3-4.
Cnn ② Pairs + 3, 4, 5, ...

- train kernels of size $K \times 100$
where $K = \text{width in terms of \# words}$.

- Recall: kernels "scan" (sweep) across picture
the field of a picture; // *review what is learned in picture cnns* →

- in this case ~~the~~ $K \times 100$ kernel would
just sweep across the sentence one.



- it would be trained to look for
a 2-word "pattern" of meaning
that would contribute to learning subjective/
objective

- is one enough? - no, suspect
would need several

- is one size enough → maybe not

→ A2 suggests have N_1 kernels of size $K_1 \times 100$
(kernel) N_2 " " " $K_2 \times 100$
- each one randomly initialized.

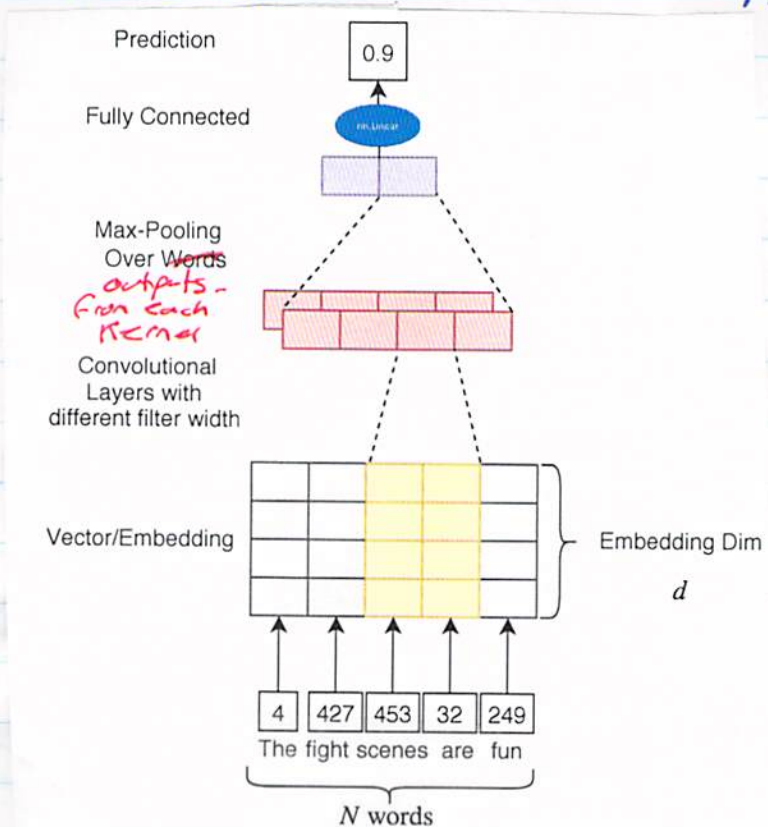
- what is the output size that you get
from an input sentence of N words and
a kernel of size 2×100 ? (where 100 is embedding
dim)
- assume that stride = 1 (recall?)

→ get $N-1$ values out. - get N_1 sets
- similar for K_2 , get $N-K_2+1$ values.

⇒ $Y_{con}[k]$ in A2 suggests picking the maximum
across all $N-K+1$ values ⇒ i.e. max pool!

⇒ feed all values from all kernels
into an MLP.

- does it make sense what is happening?)



- it should look for patterns of 1-6 words or so that give ~~input~~ sense of objective or subjective.

- how could you know what the kernels are trained to look for? - closest words:
 - discuss. *↳ ie what do they learn?*

→ asked to do this in A2: these should be textual features just like ~~vision~~ vision CNN has visual features.

Recurrent Neural Networks

→ [Aside RNNs have traditionally been used for text
 → have been replaced by Transformers.
 → were problematic all along → convergence tricky; LSTM, GRU unclear
 → create bottleneck for information flow → hidden vector.
 → will skip in favour of:

- another way to look at what the neural net is looking for is to ask which word(s) were significant in the decision using ^{the} integrated gradient method.
- talk through details of assignment
 - dataset / tokenization / iterators are given
- gradco - show / motivate

Lecture 3 Addenda.

- in introduction on embeddings.
- words are now numbers.
- just like pictures became numbers.
→ e.g. each pixel is R, G, B numbers.
- each pixel has 3-dimensions Red Green, Blue.
- each word has 200 in A1B3 800 in A154.
Glove has 50 → 100 → 300p.

⇒ far more information in 1 word's embedding than 1 pixel - e.g. 100's numbers vs 3.

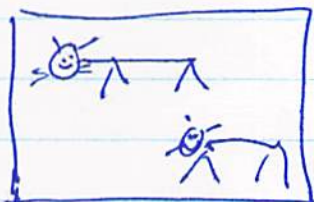
→ those 100 numbers might contain all meaning?

- What about "1 picture is worth 1000 words"
- picture has many pixels

- baseline works surprisingly well - if keep training will get better

- CNN Review

picture



- CNN has kernels that are convolved with picture.

- Kernel is learned through back prop.

- have multiple

e.g. etc