

10-315 Introduction to ML

Natural Language Processing

Instructor: Pat Virtue

Natural Language Processing (NLP)

Practical Deep Learning

NLP Intro

Feature Engineering/Learning for Text

N-gram Language Models

Word Embedding Language Models

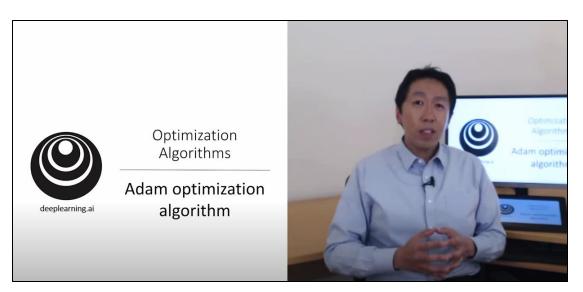


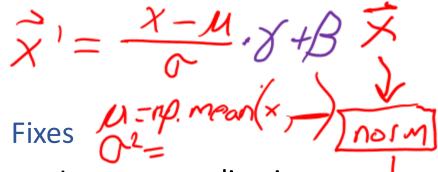
Practical Deep Learning

Issues

- Numerical Stability
- Vanishing/Exploding gradients
- Overfitting

Andrew Ng videos on Hyperparameter Tuning!





- Input normalization
- Weight initialization
- Batch normalization
- Adam
 - Exponential moving average
 - Momentum
 - RMSprop
- Learning rate decay
- Skip connections
- Dropout

Issues

- Numerical Stability
- Vanishing/Exploding gradients

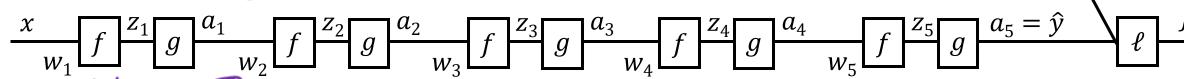
Fixes

- Input normalization
- Weight initialization
- activations = Batch normalization

 activations = Batch normalization

 feature map to autoences.

 (latent variable) = y



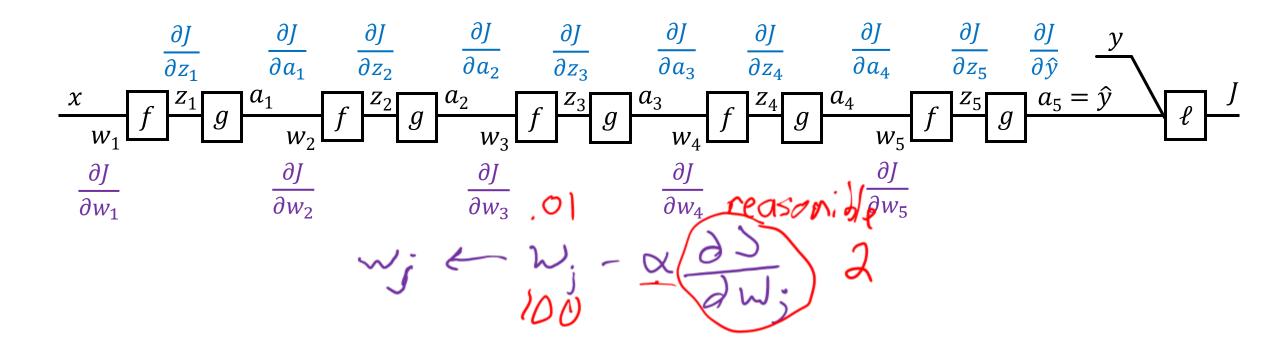


Issues

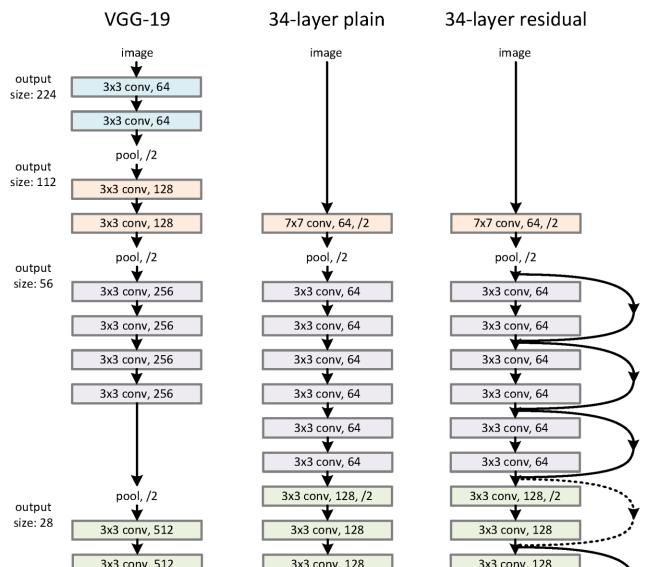
- Numerical Stability
- Vanishing/Exploding gradients

Fixes

- Input normalization
- Weight initialization
- Batch normalization



Issues



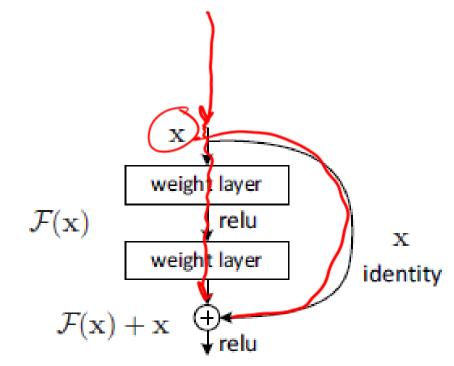


Figure 2. Residual learning: a building block

Natural Language Processing (NLP)

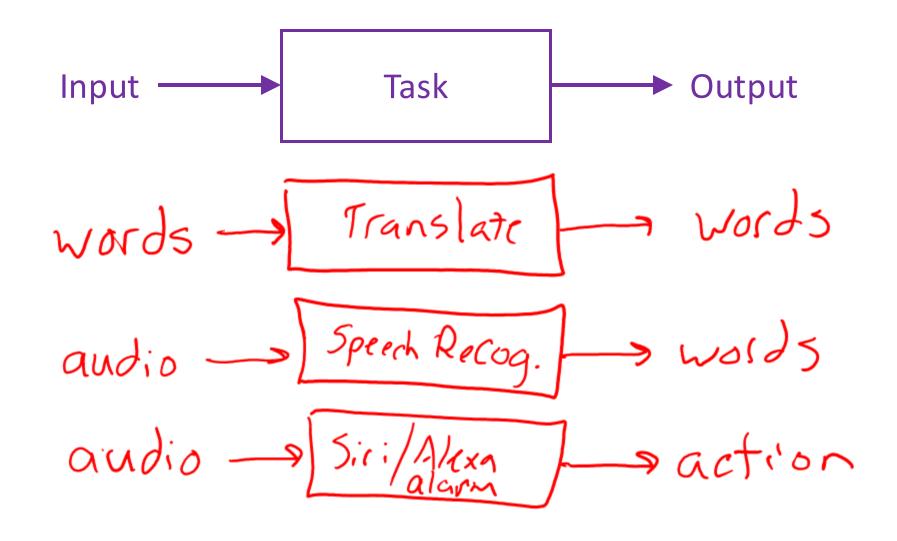
Language is Hard

Emphasis can drastically change meaning

I didn't eat your dog

NLP Tasks Examples

How many different NLP Input/Output agents can you think of?



NLP Task: Sentiment Analysis

Sentiment analysis demo

https://text2data.com/Demo

"I recommend that you find something better to do with your time"

I recommend that you find something better to do with your time

This document is: **positive (+0.60)**



Mw.

11 i vill not eat"

Text Features

SPAM Classification

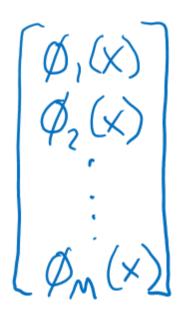
 $X = \{ \{ \{ \{ \{ \} \} \} \} \}$

SPAM Classification

Bag of words: Vector of length the size of the vocabulary

Two options

- Word occurrence: Binary: does the word exist (at least once) or not
- Word histogram: Integer: count of how many times the word appears



Predicting Rating from Written Movie Review

$$x = \frac{1}{2} \text{ text}$$

$$\phi(x) \in \mathbb{R}^{M}$$

$$\phi_{1}(x) = \# \text{ times } \text{ fantastic}$$

$$\phi_{2}(x) = \# \text{ times } \text{ fantastic}$$

$$\phi_{n}(x) = \# \text{ times } \text{ fantastic}$$

$$\phi_{n}(x) = \# \text{ times } \text{ fantastic}$$

$$\dot{y} = h(\dot{x})$$

$$= g(\dot{\partial} \dot{x})$$

$$= g(\dot{\partial} \dot{x})$$

Text encoding: token → index within vocabulary

What is our "vocabulary", i.e., what are our tokens?

- Byte pair encoding
- Words and punctuation

Corpus

am Sam. I am

Sam. Sam I am.

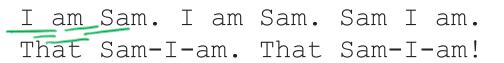
That Sam-I-am.

That Sam-I-am!

Text encoding: token → index within vocabulary

What is our "vocabulary", i.e., what are our tokens?

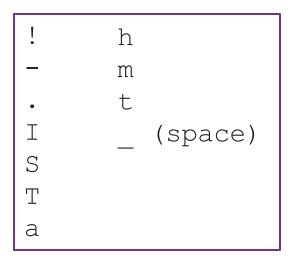
Character token BPE token Word token vocabulary vocabulary vocabulary 0:! 11: Sam 1: -1: -12: Th 1:.. 2: . 2: . 13: am 2: am 3: I 14: Sam 3: T 3: I 4: S 15: am 4: S 4: Sam 5: T 5: T 5: That 6: a 6: a 6: Sam-I-am 7: h 7: h 8: m 8: m 10: _(space) (space)



Text encoding: token → index within vocabulary

Byte pair encoding (BPE)

Current token vocabulary (Initialize to characters in corpus)



Pair frequencies

Count frequencies of all pairs of current tokens appearting together in corpus



New tokens

Add most frequent pair to tokens

```
! h
- m
. t
I
Space)
Sam
T
a
```

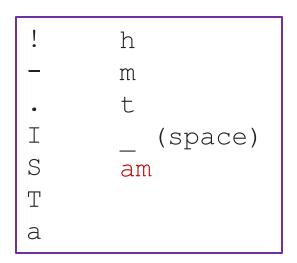


I am Sam. I am Sam. Sam I am. That Sam-I-am!

Text encoding: token → index within vocabulary

Byte pair encoding (BPE)

Current token vocabulary (Initialize to characters in corpus)



Pair frequencies

Count frequencies of all pairs of current tokens appearting together in corpus

```
_S: 4
_I: 1
Th: 2
ha: 2
at: 2
_am: 3
Sam: 5
```

New tokens

Add most frequent pair to tokens

```
! h
- m
. t
I __ (space)
S am
T Sam
a
```

Text encoding: token → index within vocabulary

What is our "vocabulary", i.e., what are our tokens?

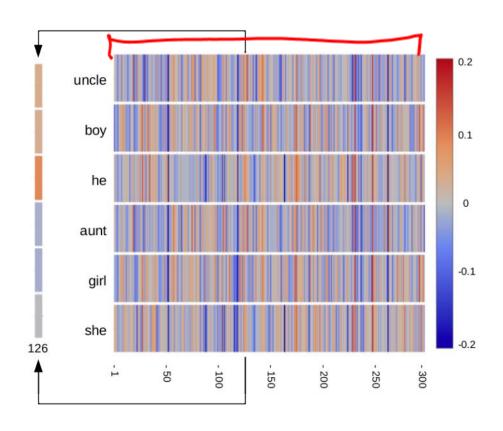
Character token BPE token Word token vocabulary vocabulary vocabulary 0:! 0:! 0:! 11: Sam 1: -1: -12: Th 1: . 2: . 2: . 13: am 2: am 3: I 3: I 14: Sam 3: I 4: S 4: S 15: am 4: Sam 5: T 5: T 5: That 6: a 6: a 6: Sam-I-am 7: h 7: h 8: m 8: m 9: t (space) 10: (space)

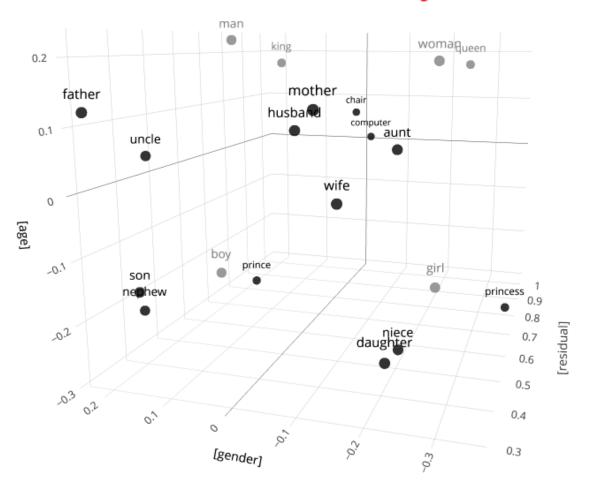
Text Feature Learning

W

Feature learning

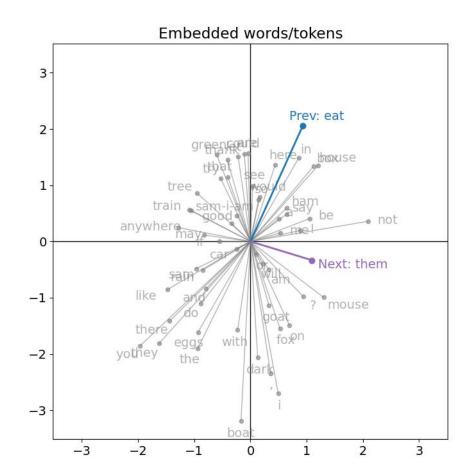
Word to Vec

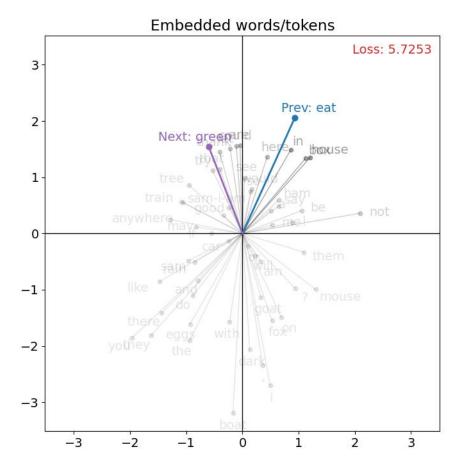




Text Feature Learning

Word embeddings





Feature Engineering vs Feature Learning

Feature engineering

- Humans decide what the features may be useful
- Humans implement feature extraction algorithms

Feature learning

- Humans choose data and performance measure
- Humans decide on structure/algorithm to learn features
 - Features are just intermediate values between input and output
 - Structure defines number of feature values
- Allow machine to map data to feature values as needed

Feature Engineering vs Feature Learning

Feature engineering

- Humans decide what the features may be useful
 - Leverage human experience for that input
 - Consider data available
- Humans implement feature extraction algorithms
 - Leverage data and compute power

Pro: Human interpretable features \rightarrow trained models easier to explain

Con: Features may not be sufficient for effective training

Con: Humans likely needed to adapt to new tasks

Pro: Less likely to overfit as humans have selected impactful features

Feature Engineering vs Feature Learning

Feature learning

- Humans choose data and performance measure
- Humans decide on structure/algorithm to learn features
 - Features are just intermediate values between input and output
 - Structure defines number of feature values
- Allow machine to map data to feature values as needed

Con: Humans cannot interpret features \rightarrow trained models unexplainable

Pro: Allow machine to search larger space for efficient features

Pro: Humans may not be needed to adapt to new tasks (just new/more data)

Con: Can overfit to data as there is no human logic in feature definition

Natural Language Processing (NLP)

Practical Deep Learning

NLP Intro

Feature Engineering/Learning for Text

N-gram Language Models

Word Embedding Language Models



Language Models

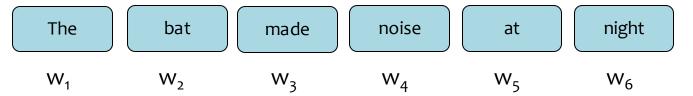
Language Models

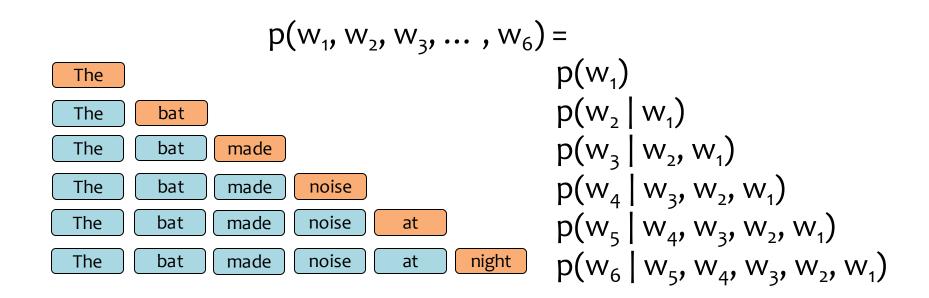
N-grams

N-gram Exercise

The Chain Rule of Probability

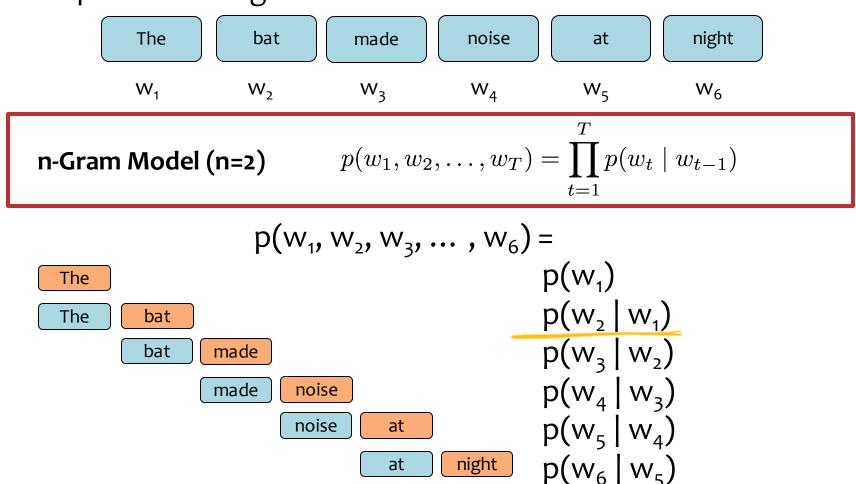
<u>Question</u>: How can we **define** a probability distribution over a sequence of length T?





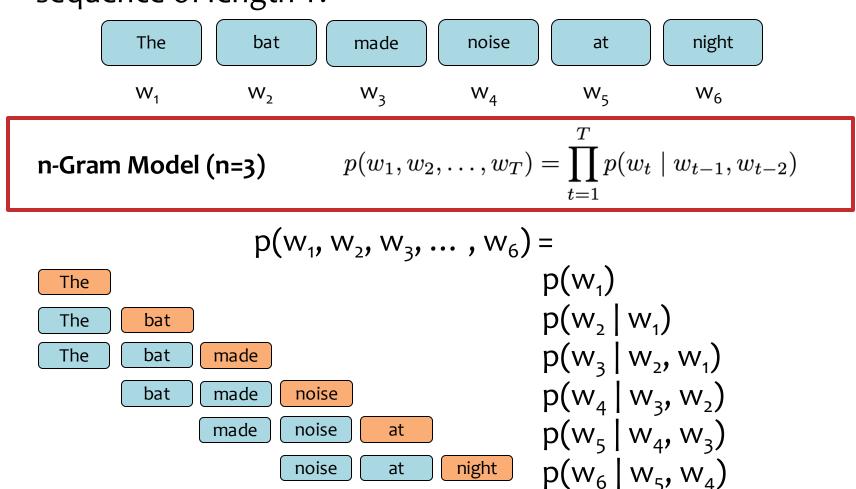
n-Gram Language Model

<u>Question</u>: How can we **define** a probability distribution over a sequence of length T?



n-Gram Language Model

<u>Question</u>: How can we **define** a probability distribution over a sequence of length T?



Where do the n-gram probabilities come from?

<u>Google n-grams demo</u>

N-gram probabilities

Vocabulary size: 30,000

Self-supervised

Example: Jane Austen, Pride and Prejudice

Vanity and pride are different things, though the words are often used synonymously. A person may be proud without being vain. Pride relates more to our opinion of ourselves, vanity to what we would have others think of us.

Self-supervised learning (auto-regressive)

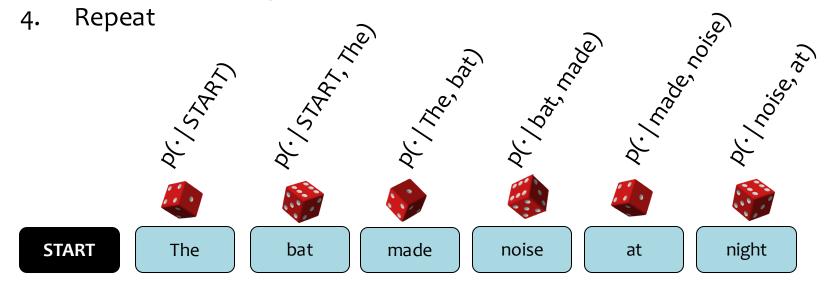
Example: Jane Austen, Pride and Prejudice

Vanity and pride are different things, though the words are often used synonymously. A person may be proud without being vain. Pride relates more to our opinion of ourselves, vanity to what we would have others think of us.

Sampling from a Language Model

<u>Question</u>: How do we sample from a Language Model? Answer:

- 1. Treat each probability distribution like a (50k-sided) weighted die
- 2. Pick the die corresponding to $p(w_t | w_{t-2}, w_{t-1})$
- 3. Roll that die and generate whichever word w_t lands face up



N-gram Examples

Random samples from language model trained on Shakespeare:

N=1: "in as , stands gods revenge ! france pitch good in fair hoist an what fair shallow-rooted , . that with wherefore it what a as your . , powers course which thee dalliance all"

n=2: "look you may i have given them to the dank here to the jaws of tune of great difference of ladies . o that did contemn what of ear is shorter time; yet seems to"

n=3: "believe , they all confess that you withhold his levied host , having brought the fatal bowels of the pope ! ' and that this distemper'd messenger of heaven , since thou deniest the gentle desdemona ,"

n=7: "so express'd : but what of that ? 'twere good you do so much for
charity . i cannot find it ; 'tis not in the bond . you , merchant , have
you any thing to say ? but little"

This is starting to look a lot like Shakespeare... because it is Shakespeare

Slide: CMU, Zico Kolter