## Sequence Learning Feed-Forward Networks for Sequence Data

### **Korbinian Riedhammer**

TECHNISCHE HOCHSCHULE NÜRNBERG GEORG SIMON OHM

# Feed-Forward Networks ....on sequence data



many-to-one



### many-to-many

## **Context is Crucial**

### Example: sentiment classification



### Solution: Use context windows to learn temporal relations



### **Connectionist HMM**



Renals et al., 1992: Connectionist Probability Estimation in the DECIPHER Speech Recognition System

• Observation: At decoding time, we need emission probabilities of all (active) states

• Problem: GMMs don't generalize well

 Idea: Use NN to "predict" emission probs for all states at the same time

### **Connectionist HMM**



- requires alignment
- "one-hot encoding"
- cross-entropy loss

↓ 
$$b_1(...)$$
 0  
 $b_2(...)$  1  
.... 1  
 $b_S(....)$  0

 $X_1$ 





- Bi-grams probabilities limit the cont
- How could we learn (not count) these?

### • Recall n-gram probabilities: count observed ngrams, use back-off for unseen

**text:** 
$$P(w_1, w_2, ..., w_n) = P(w_1) \prod_{i=2}^{N} P(w_i | w_{i-1})$$

## Why word-embeddings?

Very enjoyable nonsense, this movie



Neutral

Positiv

Negativ

## **One-hot representation**

very	enjoyable	nonsense	this	movie
1	0	0	0	0
			0	0
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1

## **One-hot representation**

very	enjoyable	nonsense	this	movie	film
1	0	0	0	0	
0	1	0	0	0	
0			0	0	
0	0	1	0	0	
0	0	0	1	0	
0	0	0	0	1	

## **One-hot representation**

Problems:

- $\rightarrow$  No relationships between words
- (e.g., synonyms like film/movie)
- → Vocabulary size explodes

very	enjoyable	nonsense	this	movie	film
1	0	0	0	0	0
0	1	0	0	0	0
0	0	1	0	0	0
0	0	0	1	0	0
0	0	0	0	1	0
0	0	0	0	0	1

## How to improve?

- fixed size vectors
- meaningful representations

do

g	movie	film

## How to improve?

- words
- meaning encoded in values
- distributed representations

dog	movie	film		
	0.9	0.8	0.8	"moves"
	0.0	0.6	0.6	art
	0.9	0.8	0.2	US-English
	0.0	0.0	1.0	creature
	1.0	1.0	0.5	noun

How would you automatically generate distributed representations?



...

Behind the tree **<u>hides</u>** a **<u>hairy</u>**, **small <u>Wolpertinger</u>**.





...

Behind the tree **hides** a **hairy**, **small Wolpertinger**.

A **small tabby** cat **hides** behind the barn.





Behind the tree **hides** a **hairy**, **small Wolpertinger**.

A small tabby cat hides behind the barn.

A <u>scruff little dog hides</u> under the car.





"You shall know a word by the company it keeps."

J.R. Firth, A synopsis of linguistic theory 1930-55, 1957



- General idea:
  - Embeddings can be automatically learnt from data
  - Enough data represents covers many relationships
  - Include the context / context words

I would like a glass of apple juice.

An apple grows on the tree.

Yesterday, my father baked an apple pie.

She drank a glass of orange juice.

There is an orange tree in the backyard.

First, peel the orange.



target word:	0
	0
movie	0
	0
	0
	0
	1



	0	
	0	
	0	
	1	$W_{C}$
••••		U
	0	
	0	
	0	
		0 0 0 1  0 0

0.1

0.2

. . .

0.9

0.0











Where do the projection matrices  $W_T$ and  $W_C$  come from?  $\rightarrow$  They have to be learned!



### Word2Vec Skip-gram

- Skip-gram:
  - choose context words to generate positive samples
  - around the target word
  - Example:
    - Let's go see a movie at the cinema

### must be in relationship to target word, e.g., environment of +/- 2 words





### Word2Vec **Negative sampling**

- Negative sampling:
  - choose random words from the vocabulary
  - label as negative samples
  - Sampling frequency depending or of words in the dataset
  - Let's go see a movie at the cinem

on the frequency	Zielwort	Kontextwort	Label
	movie	see	1
a>	movie	dear	0
	movie	autotomy	0
	movie	where	0



### Word2Vec Training





### Word2Vec Training





### Word2Vec Where will embeddings be extracted?







- independent of vocabulary size
- smaller dimensionality than vocabulary size
- representation of relationships between words



### Word2Vec Problem solved

### Very enjoyable nonsense, this movie



very	enjoyable	nonsense	this	movie
0.6	0.01	0.03	0.3	0.01
0.02	0.9	0.32	0.88	0.12
0	0.2	0.25	0	0.25
0.22	0.33	0.8	0.1	0.2
0.88	0.65	0.23	0.24	0.1
0.01	0.23	0.65	0.44	0.9



### **Attention:**

### $W_T$ is usually pre-trained on large databases, only "fine-tuning" necessary later



### Continuous Bag of Words (CBOW) Predict center word given context



Mikolov et al., 2013. "Efficient Estimation of Word Representations in Vector Space"



Visualization of semantic relationships of words;

**Good embeddings encode semantic** relationships



Male-Female

### Verb tense

### **Country-Capital**

https://www.tensorflow.org/images/linear-relationships.png



### Word2Vec Limitations

- Out-of-Vocabulary
  - Also: typos, compounds



• Also: slang, shortening



### Shared radical eat eats eaten eater eating

Figures: https://amitness.com/2020/06/fasttext-embeddings/



### FastText

- Solution:
  - Use sub-words (character n-grams) instead
  - Re-use skip-gram and negative sampling
  - Bojanowski 2017: 3-6 grams

• Observation: Words are inherently a problem (OOV, typos, morphology, etc.)

Bojanowski, Grave, Joulin and Mikolov, 2017: Enriching Word Vectors with Subword Information

### FastText **Step 1: Decompose to Sub-Words**

- eating —> <eating> • Enclose any word in the training set with <>
- Extract character n-grams with sliding window



Use hashing to reduce memory; count for bin instead of actual token 



- <eating>
- 3-grams <ea eat ati tin ing ng>



### FastText **Step 2: Modify Skip-Gram & Negative Sampling**

- Sum up the n-gram vectors and the vector of the actual word
- Sample positive and negative context (word vectors)
- Compute dot-product for actual and negative context, and use SGD to update parameters



### FastText Insights

 Impro analo rich la

			word2vec- skipgram	word2vec- cbow	fa
oves performance on <b>syntactic word</b> <b>Dav tasks</b> significantly for morphologically		Czech	52.8	55.0	77
anguage like Cz	ech and German	German	44.5	45.0	56
<u>Singular/plural</u>	Base/Comparative	English	70.1	69.9	74
cat → cats	good → better	Italian	51.5	51.8	62
uug	rough —> r	L		1	

 Degrades performance on semantic ana tasks compared to Word2Vec.



woman — > queen

		word2vec-skipgram	word2vec-cbow	fas
alogy	Czech	25.7	27.6	27
	German	66.5	66.8	62
	English	78.5	78.2	77
	Italian	52.3	54.7	52





# **FastText**Insights

- Using sub-word information with character-ngrams has better performance than CBOW and skip-gram baselines on wordsimilarity task.
- Representing out-of-vocab words by summing their subwords has better performance than assigning null vectors.

		skipgra	cbo	FT null	FT cha
Arabic	WS353	51	52	54	
	GUR35	61	62	64	•
German	GUR65	78	78	81	
	ZG222	35	38	41	
English	RW	43	43	46	
	WS353	72	73	71	•
Spanish	WS353	57	58	58	
French	RG65	70	69	75	•
Romani	WS353	48	52	51	
Russian	HJ	69	60	60	



### **Time-delay Neural Networks** Waibel et al. 1989



Peddinti et al., 2015. "A time delay neural network architecture for efficient modeling of long temporal contexts"

- Frames are typically features (MFCC, word embeddings, …)
- Concatenate frames to form contexts
- Go from narrow to wide with layers
- Lower layers learn "local" features
- Higher layers learn temporal relationships

## ConvNets



- Motivation:
  - Convolution of signal with special kernels can be a great feature
  - Well established in computer graphics (eg. Sobel edge detector)
- 1D time series: 1D convolutions
  - "within-feature convolutions"
- 2D image: 2D convolutions
  - "across-feature convolutions"

Dumoulin, V. and Visin, F. "A guide to convolution arithmetic for deep learning"





12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0



## **ConvNets Building Blocks**

- Convolution:
  - kernel size, eg. 3x3, 1x3
  - stride, step size, eg. 1
  - padding, what to do at the edges? eg. zero-pad
- Pooling to reduce/increase resolution
  - average, max, ...





### **Historic Note**

- TDNN (1989): effectively 1D convolutions
- 4.7%)



### LeCun at al., 1998: LeNet-5 architecture, MNIST error rate 0.8% (regular FF:

### **Recap** Feed-Forward Networks for Sequence Data

- Use context windows, eg. by concatenation
- Use embeddings for discrete symbols (which effectively use 1-hot)
- Use convolutions (1D, 2D) to extract temporal structure from context window
- Works for all modalities:
  - Audio: eg. MFB, MFCC
  - Text: Word Vectors