# Neural language models

Daniel Hsu COMS 6998-7 Spring 2025

# NLP (Almost) from Scratch

Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, Pavel Kuksa; JMLR 12(76):2493–2537, 2011.

# Low-level NLP benchmarks circa 2000-2010ish

- POS (parts-of-speech tagging)
- Chunking
- NER (named entity recognition)
- SRL (semantic role labeling)



• Train tagging models using supervised learning

# Architectures: window approach

Input Window			/	word (	of interest		
$\operatorname{Text}$	$\operatorname{cat}$	$\operatorname{sat}$	on	the	mat		
Feature 1	$w_1^1$	$w_2^1$			$w_N^1$		
:							
Feature K	$w_1^K$	$w_2^K$			$w_N^K$		
Lookup Tabla							
$LT_{W^1} \longrightarrow$							
:							
$LT_{W^K} \longrightarrow$							
	concat						
Linear							
$M^1 \times \odot \longrightarrow$							
	~		$n_{hu}^1$		$\rightarrow$		
HardTanh							
Linear							
$M^2 \times \stackrel{\frown}{\odot} \checkmark$							
	$\overleftarrow{n_{hu}^2 = \# \text{tags}}$						

#### Architectures: sentence approach



# What's hard-coded and what's trainable?

#### • <u>Hard-coded</u>:

- Use of discrete word features
- Extra discrete features
  - Some are generic (e.g., "stemmed" versions of words)
  - Some are task-specific (e.g., in SRL: "distance to target word")
- Architecture
  - E.g., window size, choice of non-linear activation functions
- <u>Trainable</u>:
  - Embedding vectors ("lookup tables") for discrete features
  - Weight parameters for two linear layers

# Supervised learning

- Use standard labeled data sets
- Training objectives:
  - Word-level conditional log-likelihood
  - Sentence-level conditional log-likelihood based on Viterbi-like decoding
- <u>Results</u>:
  - Far from state-of-the-art (circa 2008-2011)

### Using unlabeled data

- <u>Unlabeled data</u>: English Wikipedia; Reuters "RCV1" news articles
- <u>Goal</u>: "train language models that compute scores describing the acceptability of a piece of text"
  - Use "window approach" architecture
  - Computational burden in standard training approach:

$$\log \hat{P}(w_{t}|w_{t-n+1}, \dots, w_{t-1}) = \log \frac{\exp(f_{\theta}(w_{t}, w_{t-1}, \dots, w_{t-n+1}))}{\sum_{w} \exp(f_{\theta}(w, w_{t-1}, \dots, w_{t-n+1}))}$$

Summation over entire vocabulary!

• Contrastive ranking alternative: encourage

 $f_{\theta}(w_{t-2}, w_{t-1}, w_t, w_{t+1}, w_{t+2}) > f_{\theta}(w_{t-2}, w_{t-1}, w, w_{t+1}, w_{t+2})$  for uniformly random w

# Semi-supervised learning



Language model (trained using unlabeled data)

POS tagger (trained using labeled data)

#### Results

Approach	POS	CHUNK	NER	SRL
	(PWA)	(F1)	(F1)	(F1)
Benchmark Systems	97.24	94.29	89.31	77.92
NN+WLL	96.31	89.13	79.53	55.40
NN+SLL	96.37	90.33	81.47	70.99
NN+WLL+LM1	97.05	91.91	85.68	58.18
NN+SLL+LM1	97.10	93.65	87.58	73.84
NN+WLL+LM2	97.14	92.04	86.96	58.34
NN+SLL+LM2	97.20	93.63	88.67	74.15

LM1 = only using Wikipedia, small vocab LM2 = Wikipedia + Reuters, big vocab, longer training

# Multi-task learning



# Theory

Ando and Zhang (JMLR 2005), "A Framework for Learning Predictive Structures from Multiple Tasks and Unlabeled Data"

- Goal of semi-supervised learning:
  - Use unlabeled data to identify "good functional structures" --- i.e., ways to represent data so that target function becomes "smooth" or "simple"
- Why multi-task learning can help:
  - "When we observe multiple predictors for different problems, we have a good sample of the underlying predictor space, which can be analyzed to find the common structures shared by these predictors"
- <u>Method</u>: Train representations using "self-supervised" auxiliary tasks based on unlabeled data; use representations for supervised tasks