## Neural Language Models

#### The New Frontier of Natural Language Understanding

Gabriele Sarti

#### University of Trieste, SISSA & ItaliaNLP Lab

StaTalk 2019









## **Table of Contents**

- Natural Language Processing: On a Quest for Meaning
- Modeling Natural Language: Why's and How's
- The Challenges of True Understanding



## **A Problem of Representations**

Representation Learning is central for AI, neuroscience and semantics.



Figure 1: Hierarchy of features visualized for a CNN trained on ImageNet (Wan et al. 2013).

- For images, hierarchical representations exploiting locality of features.
- What about language? Not so easy!
- Distributional Hypothesis: Semantically related words are distributed in a similar way and occur in similar contexts.



Introduced in linguistics by <u>Harris</u> <u>1954</u>, currently explored in cognitive science.



Figure 2: Linear word relations from <u>Tensorflow Tutorials</u>.

## **Early Years: Statistical Representations for NLP**



 $w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$ 

 $tf_{ij}$  = number of occurrences of *i* in *j*  $df_i$  = number of documents containing *i* N = total number of documents

**Figure 3:** Some examples of statistical and machine learning approaches to learn text representations. *From left to right:* one-hot encoding of vocabulary terms, sentence-level lexical, syntactic and morpho-syntactic features and term frequency-inverse document frequency (tf-idf) formula.

## **Recent Times: Unsupervised Representations**

- Word embeddings: Dense vector representations of words learned by optimizing a loss function.
- Main problems: biases and disambiguating polysemy.

Well-known examples of pretrained static embeddingsare Word2Vec (Mikolov et al. 2013), GloVe (Penningtonet al. 2014) and FastText (Bojanowsky et al. 2017).



**Figure 4**: A visual representation of the Skip-gram method used to train Word2Vec embeddings. Input is a one-hot of each pair of target-context word in the sliding window. W and W' are target and context representations.

## **Context is Key for Meaning**

- Contextual Embedding: Embeddings as functions of the entire input sentence.
- Idea: Use a task to induce contextual representations inside a neural network exploiting sentence information.

Introduced by CoVe (<u>McCann et al. 2017</u>) for the machine translation task, popularized by ELMo (<u>Peters et al. 2018</u>) for language modeling.



Figure 5: A bidirectional LSTM (<u>Hochreiter and</u> <u>Schmidthuber, 1997</u>) representing the base model used for ELMo contextual word embeddings. ELMo is a task-specific combination of the internal representations in the biLSTM and uses regular and backward LM.

## The Language Modeling Task

Language Modeling (LM): Predict future token given history.

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}) \qquad \mathcal{L} = -\sum_{i=1}^n \log p(x_i | x_1, \dots, x_{i-1}; \Theta)$$

**Figure 6:** From left to right, the joint probability of a sentence defined as single word probabilities given previous context and the loss function used for LMs. Minimizing the log likelihood corresponds to maximizing the probability of correctly guessing words.

 Why LM? Unsupervised, requires knowledge and improve generalization. Alternatives: Masked Language Modeling and MT.

## LMs are Unsupervised Multitask Learners

Problem: ELMo still require task-specific models to leverage contextual embeddings.

#### Solutions:

- Task-specific fine-tuning in ULMFiT (<u>Howard & Ruder</u> <u>2018</u>) inspired by ImageNet pre-training for CV tasks.
- Generative pre-training of a transformer LM (<u>Radford</u> <u>et al. 2018</u>) with possible supervised fine-tuning.
- Results: SOTA on most language-related tasks, from sentiment analysis to NER.



**Figure 7**: The transformer decoder in OpenAI GPT. The grey block represent the transformer block described in <u>Vaswani et al. 2017</u> and can be stacked.

## **Attention Is All You Need**

- RNNs are problematic since hidden states must be computed sequentially.
- Attention mechanisms were used in conjunction with RNN to capture long-range relations, inspired by MT.
- Transformers (<u>Vaswani et al. 2017</u>) use only attention and fully connected layers to create highly scalable networks capturing distant patterns.

 $\begin{aligned} \text{Attention}(Q, K, V) &= \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V \end{aligned} \qquad \begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O \\ \text{where head}_i &= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \end{aligned}$ 

Figure 8: Scaled dot-product self-attention introduced by Vaswani et al. Queries Q and keys K have dimension  $d_k$ This type of attention is efficient thanks to matrix multiplications and can be augmented with multiple heads capture information from different representation subspaces.

## More Data and Parameters Are All You Need (?)



### What Does it Mean to Understand Language?



## **Modeling is not Understanding**

- Perplexity: Exponentiation of the entropy of a discrete probability distribution. How "unsure" the model is in predicting next event.
- Used as a quality measure for language models, the lower the better.

$$2^{H(s)} = 2^{-rac{1}{N}\sum_{i=1}^N \log_2 p(w_i)} = (2^{\sum_{i=1}^N \log_2 p(w_i)})^{-rac{1}{N}} = (p(w_1) \dots p(w_N))^{-rac{1}{N}}$$

Figure 10: Perplexity of a sentence s assuming each word has the same frequency -1/N.

- Lower perplexity doesn't imply better understanding. Performance on NLU tasks can be improved by statistical cues (<u>Niven & Kao 2019</u>)
- Need to evaluate understanding and generalization in other ways.

#### **Current Directions: NLU & NLI Benchmarks**

## GLUE SuperGLUE @decaNLP SWAG & Hella Swaq

**Figure 11:** Some of the most popular benchmarks used to evaluate LM generalization capabilities. GLUE and SuperGLUE (<u>Wang et al.</u>) focus on language understanding tasks, decaNLP (<u>McCann et al. 2018</u>) is a set of 10 general NLP tasks and SWAG/HellaSWAG (<u>Zellers et al.</u>) focus on inference with adversarial filtering and grounded dialogue.

Rank	Name	Model	URL	Score	BoolQ	СВ	COPA	MultiRC	ReCoRD	RTE	WiC	WSC	AX-g	AX-b
1	SuperGLUE Human Baselines	SuperGLUE Human Baselines		89.8	89.0	95.8/98.9	100.0	81.8/51.9	91.7/91.3	93.6	80.0	100.0	99.3/99.7	76.6
2	T5 Team - Google	Τ5		88.9	<mark>91.</mark> 0	93.0/96.4	94.8	88.2/62.3	93.3/92.5	92.5	76.1	93.8	92.7/91.9	65.6
3	Facebook Al	RoBERTa		84.6	87.1	90.5/95.2	<mark>90.6</mark>	84.4/52.5	90.6/90.0	88.2	69.9	<mark>89.0</mark>	91.0/78.1	57.9
4	IBM Research Al	BERT-mtl		73.5	84.8	89.6/94.0	73.8	73.2/30.5	74.6/74.0	84.1	66.2	<mark>61.0</mark>	97.8/57.3	29.6

**Figure 12:** SuperGLUE leaderboard as of November 18, 2019. Despite having been created to be tricky for transformer LMs, current models are approaching human performances, suggesting a new direction will soon be needed.

## **Interpreting and Probing Language Models**

- Explainability is a common trend in black-box deep learning approaches.
- For NLP models, its declinations are:
  - Probing language models for linguistic information (<u>Hewitt et al. 2019, Jawahar et al.</u> <u>2019, Lin et al. 2019, Tenney et al. 2019</u>).
  - Study attention activations and the evolution of representations (<u>Voita et al. 2019, Vig et al. 2019,</u> <u>Michel et al. 2019, Clark et al. 19</u>)



**Figure 13:** An analysis of attention heads specialized behavior for different linguistic information. Specialized heads do the heavy lifting, the rest can be pruned (<u>Voita et al. 2019</u>).

## **Perspectives: NLP Rediscovers the Human Brain**

- Availability of cerebral data from different sources (EEG, eye-tracking, fMRI).
  - Using neuroscientific techniques (e.g. RSA by <u>Kriegeskorte et al. 2008</u>) to compare brain and LM activations (<u>Abnar et al. 2019</u>, <u>Abdou et al. 2019</u>, <u>Gauthier and Levy 2019</u>)
  - Use human signals to improve model behavior (Hollenstein et al. 2019, Barrett et al. 2018)
- Target: More parsimonious models that achieve human-like, interpretable behavior.



**Figure 14**: RSA between activations in different model layers and in a human subject brain. LSTM seems to have a more similar behavior with respect to transformers. (Abnar et al. 2019)



# Thanks for

**YOUT** softmax $(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{n}})\mathbf{V}$ 

# SISSA

#### **Gabriele Sarti**





🌐 gsarti.com





#### References

#### **A Problem of Representations**

- Harris Z., "<u>Distributional structure</u>", Words, 1954
- Firth J.R., "A synopsis of linguistic theory 1930-1955", Studies in Linguistic Analysis, 1957
- Wan et al., "Regularization of Neural Networks using DropConnect", PMLR, 2013

#### **Recent Times: Unsupervised Representations**

- Mikolov et al., "Efficient Estimation of Word Representations in Vector Space", ICLR 2013
- Pennington et al., "<u>GloVe: Global Vectors for Word Representation</u>", EMNLP 2014
- Sojanowsky et al., "Enriching Word Vectors with Subword Information", TACL 2017

#### **Context is Key for Meaning**

Hochreiter & Schmidhuber, "Long Short-Term Memory", Neural Computation 1997

- McCann et al., "Learned in Translation: Contextualized Word Vectors", arXiv 2017
- Peters et al., "<u>Deep Contextualized Word Representations</u>", NAACL 2018

#### LMs Are Unsupervised Multitask Learners

- Vaswani et al., "<u>Attention is All You Need</u>", NeurIPS 2017
- Howard & Ruder, "<u>Universal Language Model Fine-tuning for Text Classification</u>", ACL 2018
- Radford et al., "Improving Language Understanding by Generative Pre-Training", Published 2018

#### More Data and Parameters Are All You Need (?)

- Hinton et al. "Distilling the Knowledge in a Neural Network", arXiv 2015
- Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", NAACL 19
- Radford et al., "Language Models are Unsupervised Multitask Learners", Published 2019 (GPT-2)
- Liu et al., "Multi-Task Deep Neural Networks for Natural Language Understanding", ACL 2019 (MT-DNN)
- Yang et al., "<u>XLNet: Generalized Autoregressive Pretraining for Language Understanding</u>", NeurIPS 2019

- Lample & Conneau, "Cross-lingual Language Model Pretraining", NeurIPS 2019 (XLM)
- Zellers et al., "<u>Defending Against Neural Fake News</u>", NeurIPS 2019 (Grover)
- Liu et al., "<u>RoBERTa: A Robustly Optimized BERT Pretraining Approach</u>", arXiv 2019
- Sahn et al., "DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter", arXiv 2019
- NVidia, "<u>MegatronLM: Training Billion+ Parameter Language Models Using GPU Model Parallelism</u>", 2019
- Raffel et al. "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer", arXiv 19 (T5)

#### Modeling is Not Understanding

Niven & Kao, "Probing Neural Network Comprehension of Natural Language Arguments", ACL 2019

#### **Current Directions: NLU & NLI Benchmarks**

- Wang et al., "GLUE: A Multitask Benchmark and Analysis Platform for NLU", ICLR 2019
  - > "<u>SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems</u>", arXiv 2019
- Zellers et al., "<u>SWAG: A Large-Scale Adversarial Dataset for Grounded Commonsense Inference</u>", EMNLP 2018

- > "<u>HellaSwag: Can a Machine Really Finish Your Sentence?</u>", ACL 2019
- McCann et al., "<u>The Natural Language Decathlon: Multitask Learning as Question Answering</u>", arXiv 2018

#### **Interpreting and Probing Language Models**

- Hewitt et al., "<u>A Structural Probe for Finding Syntax in Word Representations</u>", NAACL 2019
- ✤ Jawahar et al., "<u>What Does BERT Learn about the Structure of Language?</u>", ACL 2019
- Lin et al., "Open Sesame: Getting Inside BERT's Linguistic Knowledge", ACL 2019
- Tenney et al., "BERT Rediscovers the Classical NLP Pipeline", ACL 2019
- Voita et al., "<u>Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned</u>", ACL 2019
- Voita et al., "<u>The Bottom-up Evolution of Representations in the Transformer: A Study with Machine</u> <u>Translation and Language Modeling Objectives</u>", IJCNLP 2019
- Michel et al., "<u>Are Sixteen Heads Really Better than One?</u>", NeurIPS 2019
- Vig et al., "<u>Analyzing the Structure of Attention in a Transformer Language Model</u>", ACL 2019
- Clark et al. "What Does BERT Look at? An Analysis of BERT's Attention", BlackBoxNLP 2019

#### **Perspectives: NLP Rediscovers the Human Brain**

- Kriegeskorte et al., "<u>Representational Similarity Analysis</u>", Frontiers 2008
- Abnar et al., "<u>Blackbox meets blackbox: Representational Similarity and Stability Analysis of Neural</u> <u>Language Models and Brains</u>", BlackBoxNLP 2019
- Abdou et al., "Higher-order Comparisons of Sentence Encoder Representations", EMNLP 2019
- Gauthier & Levy, "Linking artificial and human neural representations of language", IJCNLP 2019
- Hollenstein et al., "<u>Advancing NLP with Cognitive Language Processing Signals</u>", arXiv 2019
- Sarrett et al., "Sequence Classification with Human Attention", CoNLL 2018

#### General Inspiration and Structure of the Talk

- Weng L., "Learning Word Embeddings", Blog Post, 2017
- Weng, L., "<u>Attention? Attention!</u>", Blog post, 2018
- Weng L., "<u>Generalized Language Models</u>", Blog Post, 2019