Language Models

Andrea Esuli

Statistical Language Model

A statistical language model is a *probability distribution P* over sequences of terms.

Given a document *d* that is composed of a sequence of words $w_1 w_2 w_3$, we can define:

$$P(d) = P(w_1 w_2 w_3) = P(w_1) P(w_2 | w_1) P(w_3 | w_1 w_2)$$

Depending on the assumptions we make on the probability distribution, we can create statistical model of different complexity.

The formula above makes no assumptions and can exactly model any language, yet it is impractical because it requires to learn the probability of any sequence in the language.

Unigram model

A unigram model assumes a *statistical independence* between words, i.e., the probability of *d* is the product of the probabilities of its words:

$$P(d) = P(w_1 w_2 w_3) = P(w_1) P(w_2 | w_1) P(w_3 | w_1 w_2)$$

$$= P(w_1)P(w_2)P(w_3) = \prod_i P(w_i)$$

The bayesian classifier that uses this model is called *naïve* for this reason. Usually the models use the logs of the probabilities to work in a linear space:

$$log(\Pi_i P(w_i)) = \sum_i log(P(w_i))$$

Smoothing, e.g., add one to all frequencies, is used to avoid zero probabilities.

Bigram model

A bigram model assumes a *statistical dependence* of a word from the preceding one:

$$P(d) = P(w_1 w_2 w_3) = P(w_1) P(w_2 | w_1) P(w_3 | w_1 w_2)$$

$$= P(w_1)P(w_2 | w_1)P(w_3 | w_2) = \prod_i P(w_i | w_{i-1})$$

This simple addition is already able to capture a good amount of language regularities.

In general, the longer the n-gram we adopt for the model:

- the more semantic is captured;
- the less statistical significant is the model (memorization/generalization).

Vector Space Model

The *Vector Space Model* (VSM) is a typical machine-processable representation adopted for text.

Each vector positions a document into an *n*-dimensional space, on which learning algorithms operate to build their models

$$v(d_1) = [w_1, w_2 \dots, w_{n-1}, w_n]$$



Vector Space Model

After text processing, tokenization... a document is usually represented as vector in $R^{|F|}$, where F is the set of all the distinct *features* observed in documents.

Each feature is mapped to a distinct dimension in $R^{|F|}$ using a *one-hot* vector:

$$v('played') = [1, 0, 0, ..., 0, 0, ..., 0, 0, 0]$$

$$v('game') = [0, 1, 0, ..., 0, 0, ..., 0, 0, 0]$$

$$v('match') = [0, 0, 1, ..., 0, 0, ..., 0, 0, 0]$$

$$\vdots$$

$$v('trumpet') = [0, 0, 0, ..., 0, 1, ..., 0, 0, 0]$$

$$\vdots$$

$$v('bwoah') = [0, 0, 0, ..., 0, 0, ..., 0, 0, ...]$$

Vector Space Model

A document is represented as the weighted sum of its features vectors:

$$v(d) = \sum_{f \in d} w_{fd} v(f)$$

For example:

$$d$$
 = 'you played a good game'
 $v(d) = [0, w_{played,d}, w_{game,d}, 0, \dots \dots 0, w_{good,d}, 0 \dots \dots 0, 0]$

The resulting document vectors are sparse:

$$|\{i|v_i(d) \neq 0\}| \ll n$$

Sparse representations

$$d_1$$
 = 'you played a game'
 d_2 = 'you played a match'
 d_3 = 'you played a trumpet'

$$\begin{array}{l} v(d_1) = [0, \ w_{\text{played,d1}}, \ w_{\text{game,d1}}, \ 0 & , \ 0, \ \dots & , \ 0, \ 0 & , \ 0] \\ v(d_2) = [0, \ w_{\text{played,d2}}, \ 0 & , \ w_{\text{match,d2}}, \ 0, \ \dots & , \ 0, \ 0 & , \ 0] \\ v(d_3) = [0, \ w_{\text{played,d3}}, \ 0 & , \ 0 & , \ 0 & , \ 0, \ \dots & , \ 0, \ w_{\text{trumpet,d3}}, \ 0] \end{array}$$

Semantic similarity between features (game~match) is not captured:

 $sim(v(d_1), v(d_2)) \sim sim(v(d_1), v(d_3)) \sim sim(v(d_2), v(d_3))$

Modeling word similarity

How do we model that *game* and *match* are related terms and *trumpet* is not?

Using linguistic resources: it requires a lot of human work to build them.

Observation: co-occurring words are semantically related.

Pisa	is a	province	of	Tuscany
Red	is a	color	of the	rainbow
Wheels	are a	component	of the	bicycle
*Red	is a	province	of the	bicycle

We can exploit this propriety of language, e.g., following the *distributional hypothesis*.

Distributional hypothesis

"You shall know a word by the company it keeps!" <u>Firth (1957</u>)

Distributional hypothesis: the meaning of a word is determined by the contexts in which it is used.

Yesterday we had **bwoah** at the restaurant. I remember my mother cooking me **bwoah** for lunch. I don't like **bwoah**, it's too sweet for my taste. I like to dunk a piece **bwoah** in my morning coffee.

Word-Context matrix

A word-context (or word-word) matrix is a $|F| \cdot |F|$ matrix *X* that **counts** the frequencies of co-occurrence of words in a collection of contexts (i.e, text spans of a given length).

You cook the cake twenty minutes in the oven at 220 C. I eat my steak rare. I'll throw the steak if you cook it too much. The engine broke due to stress. I broke a tire hitting a curb, I changed the tire.

 $Context_{-2,+2}('cake') = \{['cook', 'the', 'twenty', 'minutes']\}$ $Context_{-2,+2}('tire') = \{['broke', 'a', 'hitting', 'a'], ['changed', 'the']\}$

Word-Context matrix

		Context words						
			cook	eat		changed	broke	
Words	cake		10	20		0	0	
	steak		12	22		0	0	
	bwoah		7	10		0	0	
	engine		0	0		3	10	
	tire		0	0		10	1	

Words = Context words Rows of X capture similarity yet X is still high dimensional and sparse.

Dense representations

We can learn a projection of feature vectors v(f) into a *low dimensional* space R^k , $k \ll |F|$, of *continuous space word representations* (i.e. word embeddings).

Embed: $R^{|F|} \rightarrow R^k$

Embed(v(f)) = e(f)

We force features to share dimensions on a reduced dense space Let's group/align/project them by their syntactic/semantic similarities!

SVD

Singular Value Decomposition is a decomposition method of a matrix X of size $m \cdot n$ into three matrices $U\Sigma V^*$, where:

U is an orthonormal matrix of size $m \cdot n$

 Σ is a diagonal matrix of size $n \cdot n$, with values $\sigma_1, \sigma_2 \dots \sigma_n$

V is an orthonormal matrix of size $n \cdot n$, V^* is its conjugate transpose

 $\sigma_1, \sigma_2 \dots \sigma_n$ of Σ are the singular values of X, sorted by decreasing magnitude. Keeping the top k values is a least-square approximation of X

Rows of U_k of size $m \cdot k$ are the dense representations of the features

SVD





By Nicoguaro (Own work) [CC BY 4.0], via Wikimedia Commons

GloVe

<u>GloVe: Global Vectors for Word Representation</u> is a count-based model that implicitly factorizes the word-context matrix based on the observation that the ratio of conditional probabilities better captures the semantic relations between words.

Probability and Ratio	k = solid	k = gas	k = water	k = fashion
P(k ice)	$1.9 imes 10^{-4}$	6.6×10^{-5}	$3.0 imes 10^{-3}$	1.7×10^{-5}
P(k steam)	$2.2 imes 10^{-5}$	7.8×10^{-4}	$2.2 imes 10^{-3}$	1.8×10^{-5}
P(k ice)/P(k steam)	8.9	8.5×10^{-2}	1.36	0.96

$$J = \sum_{i,j=1}^{V} f\left(X_{ij}\right) \left(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij}\right)^2 \quad \leftarrow \text{ eqn 8 of GloVe paper}$$

GloVe



Weighting function to filter out rare co-occurrences and to avoid frequent ones to dominate



embedding vectors (compressed)

co-occurrence matrix (sparse)

This part implements, as a least square problem, the equation that defines the model:

$$F(w_i, w_j, \tilde{w_k}) \approx \frac{P_{ij}}{P_{jk}}$$

Word2Vec

<u>Skip-gram and CBoW models of Word2Vec</u> define tasks of **predicting** a *context from a word* (Skip-gram) or a *word from its context* (CBoW).

They are both implemented as a two-layers linear **neural network** in which input and output words one-hot representations which are encoded/decoded into/from a dense representation of smaller dimensionality.



Word2Vec

Embeddings are a byproduct of the word prediction task.

Even though it is a prediction tasks, the network can be trained on any text, no need for human-labeled data!

The context window size ranges between two and five words before and after the central word.

Longer windows capture more semantic, less syntax.

A typical size for h is 200~300.



Skip-gram

w vectors are high dimensional, |F|

h is low dimensional (it is the size of the embedding space)

 W_I matrix is $|F| \cdot |h|$. It encodes a word into a hidden representation. Each row of W_I defines the embedding of the a word.

 W_o matrix is $|h| \cdot |F|$. It defines the embeddings of words when they appears in contexts.



Skip-gram

 $h = w_t W_I \leftarrow h$ is the embedding of word w_t

$$u = h W_o \leftarrow u_i$$
 is the similarity of h with
context embedding of w_i in W_c

Softmax converts *u* to a probability distribution *y*:

$$y_i = exp(u_i) / \sum_{j \in F} exp(u_j)$$



Skip-gram

Loss:
$$-\log p(w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2} | w_t) =$$
$$= -\log \prod_{c \in C} \exp(u_c) / \sum_{j \in F} \exp(u_j) =$$
$$= -\sum_{c \in C} \exp(u_c) + C \log \sum_{j \in F} \exp(u_j)$$

i.e., maximize probability of context

- $\sum_{c \in C} exp(u_c)$

and minimize probability of the rest

+
$$C \log \sum_{j \in F} exp(u_j)$$



Negative sampling

The $log \sum_{j \in F} exp(u_j)$ factor has a lots of terms and it is costly to compute.

Solution: compute it only on a small sample of negative examples, i.e.,

 $log \sum_{j \in E} exp(u_j)$

where words in *E* are just a few (e.g., 5) and they are sampled using a biased unigram distribution computed on training data:

$$p(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j \in F|} f(w_j)^{3/4}}$$



CBoW

CBoW stands for Continuous Bag of Word.

It's a mirror formulation of the skip-gram model, as context words are used to predict a target word.

h is the average of the embedding for the input context words.

 u_i is the similarity of h with the word embedding w_t in W_o



Which model to choose?

<u>Levy and Goldberg</u> proved that Word2Vec version that uses skip-gram with negative sampling (SGNS) implicitly computes a factorization of a variant of *X*.

<u>Levy, Goldberg and Dagan</u> ran an extensive comparison of SVD, CBoW, SGNS, GloVe.

- Results indicate no clear overall winner.
- Parameters play a relevant role in the outcome of each method.
- Both SVD and W2V performed well on most tasks, never underperforming significantly.
- W2V is suggested to be a good baseline, given its lower computational cost in time and memory.

Computing embeddings

The training cost of Word2Vec is linear in the size of the input.

The training algorithm works well in parallel, given the sparsity of words in contexts and the use of negative sampling. The probability of concurrent update of the same values by two processes is minimal \rightarrow let's ignore it when it happens (a.k.a., <u>asynchronous stochastic gradient descent</u>).

Can be halted/restarted at any time.

The model can be updated with any data (concept drift/ domain adaptation).

Computing embeddings

Gensim provides an efficient and detailed implementation.

```
sentences = [['this','is','a','sentence'],
        ['this','is','another','sentence']]
```

```
from gensim.models import Word2Vec
model = Word2Vec(sentences)
```

This is a clean implementation of skip-grams using pytorch.

Which embeddings?

Both W_1 and W_0 define embeddings, which one to use?

- Usually just W₁ is used.
- Average pairs of vectors from W_1 and W_0 into a single one.
- Append one embedding vector after the other, doubling the length.

Testing embeddings

Testing if embeddings capture syntactic/semantic properties.



Analogy test:

Paris stands to France as Rome stands to ? book ? Writer stands to painter stands to as mouse stands to Cat stands to ? cats as

e('France') - e('Paris') + e('Rome') ~ e('Italy')

a : b = c : c

$$d = \arg \max_{x} \frac{(e(b)-e(a)+e(c))^{T}e(x)}{||e(b)-e(a)+e(c)||}$$

The impact of training data

The source on which a model is trained determines what semantic is captured.

WIKI	BOOKS	WIKI	BOOKS	
se	ega	chianti		
dreamcast	motosega	radda	merlot	
genesis	seghe	gaiole	lambrusco	
megadrive	seghetto	montespertoli	grignolino	
snes	trapano	carmignano	sangiovese	
nintendo	smerigliatrice	greve	vermentino	
sonic	segare	castellina	sauvignon	

FastText word representation

FastText extends the W2V embedding model to <u>ngrams of the words</u>.

The word "goodbye" is also represented with a set of ngrams:

"<go" (star of word), "goo", "ood", "odb", "dby", "bye", "ye>" (end of word)

The length of the ngram is a parameter.

Typically all ngrams of length from 3 to 6 are included.

FastText word representation

The embedding of a word is determined as the sum of the embedding of the word and of the embedding of its ngrams.

Subword information allows to give an embedding to OOV words.

Subword information improves the quality of misspelled words.

<u>Pretrained embeddings for</u> <u>200+ languages.</u> Query word? accomodation sunnhordland 0.775057 accomodations 0.769206 administrational 0.753011 laponian 0.752274 ammenities 0.750805 dachas 0.75026 vuosaari 0.74172 hostelling 0.739995 greenbelts 0.733975 asserbo 0.732465 Query word? gearshift gearing 0.790762 flywheels 0.779804 flywheel 0.777859 gears 0.776133 driveshafts 0.756345 driveshaft 0.755679 daisywheel 0.749998 wheelsets 0.748578 epicycles 0.744268 gearboxes 0.73986

Query word? accomodation accomodations 0.96342 accommodation 0.942124 accommodations 0.915427 accommodative 0.847751 accommodating 0.794353 accomodated 0.740381 amenities 0.729746 catering 0.725975 accomodate 0.703177 hospitality 0.701426

Misspelling Oblivious Embeddings

<u>Misspelling Oblivious Embeddings</u> (MOE) are an extension of FastText embedding which explicitly model misspellings as a component of the language model.



bir + ird

bwr + wrd

A probabilistic model of misspelling is generated by observing query corrections from a query log (i.e., pairs of queries close in time with small edit distance).

The misspelling model is used to inject misspellings into text.

Multilingual Embeddings

<u>MUSE</u> (Multilingual Unsupervised and Supervised Embeddings) aligns language models for different languages using two distinct approaches:

- *supervised*: using a *bilingual dictionary* (or same string words) to transform one space into the other, so that a word in one language is projected to the position of its translation in the other language.
- *unsupervised*: using *adversarial learning* to find a space transformation that matched the distribution of vectors in the two space (without looking at the actual words).



Exploring embeddings





Word embeddings to documents

How to represent *a document* using word embeddings?

- average
- max
- max+min (double length)
- Doc2Vec
- As a layer in a more complex neural network
Doc2Vec

Proposed by <u>Le and Mikolov</u>, Doc2Vec extends Word2Vec by adding input dimensions for identifiers of documents.

```
W_I matrix is (|D|+|F|)\cdot |h|.
```

Documents ids are projected in the same space of words.

The trained model can be used to infer document embeddings for previously unseen documents - by passing the words composing them.



Exploring embeddings

Documents embedding can be used as vectorial representations of documents in any task.

When the document id is associated to more than one actual document (e.g., id of a product with multiple reviews), Doc2Vec is a great tool to model similarity between objects with multiple descriptions.



Embeddings in neural networks

An embedding layer in neural networks is typically the first layer of the network.

It consists of a matrix W of size $|F| \cdot n$, where n is the size of the embedding space.

It maps words to dense representations.

It can be initialized with random weights or *pretrained* embeddings.

During learning weights can be kept fixed (it makes sense only when using pretrained weights) or updated, to adapt embeddings to the task.

Embeddings in neural networks



Example:

https://github.com/fchollet/keras/blob/master/examples/imdb_cnn.py

Another example: <u>https://machinelearningmastery.com/use-word-embedding-layers-deep-learning-keras/</u>

Embeddings in NN: OOV words and padding

LMs that use a vocabulary do not model out-of-vocabulary words.

Add a special *unknown* word (and embedding) for such words, to be learned during training.

NNs usually process examples in *batches*, i.e., set of k examples.

Input sentences in a batch are usually required to be of the same length.

For this reason a special *padding* word (and embedding) is added before/after (be consistent!) words of shorter sentence to match the length of the longest one.

Recurrent Neural Networks

A Recurrent Neural Network (RNN) is a neural network in which connections between units form a directed cycle.

Cycles allow the network to have a memory of previous inputs, combining it with current input.

RNNs are fit to process sequences, such as text.

Text can be seen as a sequence of values at many different levels: characters, words, phrases...

$x \xrightarrow{U} s \xrightarrow{V} o$



Suggested read

Char-level LM & text generation

RNNs are key tools in the implementation of many NLP applications, e.g., machine translation, summarization, or image captioning.

A RNN can be used to learn a language model that predicts the next character from the sequence of previous ones.

The typical RNN node that is used is an Long Short Term Memory (LSTM), which is <u>robust to</u> <u>typical issues of RNNs</u>.



ELMo

<u>ELMo</u> (Embeddings from Language Models) exploits the *hidden states* of a *deep, bi-directional, character-level* LSTM to give *contextual* representations to words in a sentence.

- Contextual: representation for each word depends on the entire context in which it is used.
- Deep: word representations combine all layers of a deep pre-trained neural network.
- Character based: ELMo representations are character based, allowing the network to use morphological clues

ELMo

<u>ELMo</u> (Embeddings from Language Models) exploits the *hidden states* of a *deep, bi-directional, character-level* LSTM to give *contextual* representations to words in a sentence.



The embedding for a word is a *task-specific* weighted sum of the concatenation of the f,b vectors for each level of the LSTM, e.g.:

$$w(token) = w [f, b]_{token} + w' [f', b']_{token}$$

Context Vectors (CoVe)

<u>CoVE</u> (Context Vector), trains a sequence-to-sequence network on a Machine Translation problem (a).

The encoder of the seq-2-seq network is a two-layer bidirectional LSTM, which takes in input glove embeddings.



The encoder is then used to enrich the glove embeddings (b).

Attention-based models

<u>Attention</u> is a simple mechanism that relates the elements of two sequences so as to identify correlations among them.

Attention proved to be effective in sequence-to-sequence <u>machine</u> <u>translation models</u>.

Attention captures the contribution of each input token to determine the output tokens.



The Transformer is a network architecture for sequence-to-sequence problems, e.g., machine translation.

It replaces the traditionally used LSTM elements with a set of encoder/decoder elements based on the attention mechanism.

The transformer was proposed in the paper <u>Attention is All You Need</u>.

The elements of the transformer are at the base of many recently proposed language models.

Picture in the following slides are taken from <u>The Illustrated Transformer</u> by Jay Alammar.

The Transformer is a network architecture for sequence-to-sequence problems, e.g., machine translation.



Output Probabilities

Softmax

Linear

Add & Norm Feed

Forward

Add & Norm

Add & Norm

It replaces the traditionally used LSTM elements with a set of encoder/decoder elements based on the attention mechanism.



The encoder and decoder elements exploit attention (self-attention) to combine and transform the input from lower layers (and the last encoder in the case of the decoder).



Transformer - encoder

The encoder elements exploit self-attention to combine and transform the input from lower layers.



Self-attention correlates all the inputs between themselves.



Self-attention first transforms the embedding vector of each token into three vectors: a query vector, a key vector and a value vector.

The query represents the tokenInputas the input.EmbedThe key is used to match the otherQuerietokens against the query token.QuerieThe value is the contribution of everyKeys

The three transformations are defined by three dedicated matrices W^Q, W^K, W^V.

Input	Thinking	Machines		
Embedding	X1	X ₂		
Queries	q1	q 2		Wa
Keys	k 1	k2	100 100 100 100 100 100 100 100 100 100	Wĸ
Values	V1	V2		Wv

The query vector for a token is multiplied (dot product) with every key vector.

The factor eight stabilizes the computation.

The resulting numbers are normalized with a softmax and they become the weights for a weighted sum of the value vectors.

This concludes *single headed* self-attention.

Input	Thinking	Machines							
Embedding	X1	X2							
Queries	q 1	q ₂							
Keys	k 1	k2							
Values	V1	V2							
Score	$q_1 \cdot k_1 = 112$	q ₁ • k ₂ = 96							
Divide by 8 ($\sqrt{d_k}$)	14	12							
Softmax	0.88	0.12							
Softmax X Value	V1	V2							
Sum	Z 1	Z ₂							

Self attention can be expressed and efficiently implemented with matrices





Different sets of matrices W^Q, W^K, W^V, can be used to capture many different relations among tokens, in a *multi-headed* attention model.



Outputs from each head of the attention model are combined back into a single matrix with a vector for each token.



Z ₀ Z		Z 2			Z 3			Z 4			Z 5			Z 6			Z 7			

2) Multiply with a weight matrix W^o that was trained jointly with the model

Х

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN





Self-attention - complete view

1) This is our input sentence*

2) We embed each word* 3) Split into 8 heads. We multiply X or R with weight matrices 4) Calculate attention using the resulting Q/K/V matrices 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer





* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

R







...

...



Ζ

Residual connections

A <u>residual connection</u> simply consists in skip connections that sum the input of a layer (or more layers) to its output.

Residual connections let the network learn faster, and be more robust to the problem of <u>vanishing gradients</u>, allowing much deeper network to be used.

Residual connections are put around any attention or feed forward layer.

<u>Layer normalization</u> supports the numerical stability of the process.



Word order

In order to take into account of the actual position of tokens in the sequence, the input embeddings are added with a *positional encoding* vector.

The positional encoding function is a <u>linear transformation</u> of the token representation

$$f(t,i) = egin{cases} \sin\left(rac{t}{10000^{rac{i}{d}}}
ight), & ext{if } i \equiv_2 0 \ \cos\left(rac{t}{10000^{rac{i-1}{d}}}
ight), & ext{if } i \equiv_2 1 \end{cases}$$



Transformer - decoder

The decoder elements exploit attention and self-attention to combine and transform the input from lower layers and the last encoder.

The decoder is applied repeatedly to form the output one token at a time:



OUTPUT

Decoding time step: 1 2 3 4 5 6

Transformer - decoder

The decoder elements exploit attention and self-attention to combine and transform the input from lower layers and the last encoder.

The decoder is applied repeatedly to form the output one token at a time:



Transformer – decoder

The last linear+softmax layer converts the output of the decoder stack into probabilities over the vocabulary.

Once trainined, this last element can be replaced with other structures depending on the tasks, e.g., classification, see GPT.

Which word in our vocabulary is associated with this index?

> Get the index of the cell with the highest value (argmax)



am

5

Transformer - complete view



Transformer-based language models

Many language models stemmed from the Transformer



GPT

<u>GPT</u> (<u>Generative Pre-Training</u>), learns a language model using a variant of the Transformer's decoders that uses *masked self-attention*.



<u>GPT</u> (<u>Generative Pre-Training</u>), learns a language model using a variant of the Transformer's decoders that uses *masked self-attention*.

Masked self-attention allows the prediction to be based only on the left context of the predicted token.

GPT

GPT is trained on language generation in an unsupervised fashion, then it can be attached to additional layers to perform a supervised task.

GPT-2

<u>GPT-2</u> (paper) is a minor variant of the GPT, its main characteristic being the scale, with models up to 48 layers and 1.5 billion parameters.

It exhibited excellent capabilities of controlled text generation, opening a disucission about possible <u>unethical or illegal uses</u>.

345M Parameters

762M Parameters

1,542M Parameters

BERT

<u>BERT</u> (Bidirectional Encoder Representations from Transformers), uses the Transformer encoders as building blocks.

The model is bidirectional in the sense that the attention model can peek at both left and right contexts.

BERT

The model is pre-trained on masked word prediction and next sentence classification tasks. Fine-tuning on final task is relatively quick.

BERT

BERT produces contextualized embeddings at each level, if we use them without fine tuning, which should be used?



BERT

BERT produces contextualized embeddings at each level, if we use them without fine tuning, which should be used?

What is the best contextualized embedding for "Help" in that context?

Four Hidden

For named-entity recognition task CoNLL-2003 NER Dev F1 Score Embedding First Layer 91.0 ... Last Hidden Layer 94.9 + + Sum All 12 95.5 + Layers 5 Second-to-Last 95.6 Hidden Layer 3 12 + 11 + 10 + 9 = 2 Sum Last Four 1 95.9 Hidden Help **Concat Last**

96.1

Differences between BERT, GPT and ELMo

(from BERT paper)



Figure 3: Differences in pre-training model architectures. BERT uses a bidirectional Transformer. OpenAI GPT uses a left-to-right Transformer. ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTMs to generate features for downstream tasks. Among the three, only BERT representations are jointly conditioned on both left and right context in all layers. In addition to the architecture differences, BERT and OpenAI GPT are fine-tuning approaches, while ELMo is a feature-based approach.

Megatron-LM

Large language models require efficient distributed GPU computation.

NVidia published <u>Megatron-LM</u>, the code to efficiently train BERT_{Large} and a huge GPT2 language model:

- the BERT Large can be trained on 64 Tesla V100 GPUs in 3 days (originally 4 days on 16 TPU = 64 TPU chips, note that <u>rough estimates were negative</u> <u>about this</u>)
- the GPT2 Language model has 72 layers and 8.3 billion parameter. It is trained using a 8-way model and 64-way data parallelism across 512 GPUs.

Transformer-XL

The <u>Transformer-XL</u> model extended the Transformer model (decoder-based) in order to overcome its limitation on context length.

The key idea is use recurrent contexts during training, with fixed gradients for the tokens in the past beyond the context length.

This model work improve efficacy when working with long documents.



Other models

<u>XLNet</u>: permutation language model to overcome BERT issues:

- No [mask] during fine-tuning.
- Predictions are made independently one from the other.

"I saw a [mask] on a [mask]" "I saw a cat on a couch" "I saw a car on a road" *"I saw a car on a couch"

ERNIE: from Baidu. Continuous multitask training with task embeddings.

Software

<u>Transformers</u> provides general-purpose architectures (BERT, GPT-2, RoBERTa, XLM, DistilBert, XLNet...) for Natural Language Understanding (NLU) and Natural Language Generation (NLG) with over 32+ pretrained models in 100+ languages and deep interoperability between TensorFlow 2.0 and PyTorch.

<u>Keras Transformers</u> implements BERT and other elements of the Transformer to support the definition of Transformer-based language models.

<u>Keras-GPT-2</u> is an example of how to use pretrained GPT-2 model weights to make predictions and generate text.

BERT repo.

XLNet repo.

Summary

Language models are a powerful tool to collect unsupervised information from a language or domain and to model it in machine-readable latent representations.

Such representations allow us to:

- explore a language, a domain, and discover its relevant entities and their syntactic and semantic relations.
- infuse general knowledge about a language or a domain into a supervised learning task, enabling the learned model to express a better generalization ability.