Advanced Methods in Natural Language Processing

Session 3: Word Embeddings

Arnault Gombert

April 2025

Barcelona School of Economics

Introduction

Today's Focus: Unveiling the Power of Word Embeddings

- Sparse Vectors and Ontologies: Evolution of word representation.
- Embedding quality: Evaluating embeddings methods.

Old Fashion Embedding Techniques

- Word2Vec Insights: Understanding the mechanics and impact of Word2Vec and the Skip-Gram model.
- Exploring Static Word Embeddings: Other famous static wordembedding methods: GloVe and FastText.

Advancing to Sophisticated Embedding Techniques

- Contextual Word Embeddings: Delving into ELMo.
- Future of Embeddings: Advanced models like BERT and GPT-x.

Word Representation

Traditional word representation methods (cf. Session 1):

- One-Hot Vectors: Each token is associated with a unique index.
- Token Counts (Hans P. Luhn, 1957): Token frequencies in the text.
- **TF-IDF** (Spärck Jones, K., 1972): Token's importance relative to its frequency across documents.

Limitations of Traditional Techniques:

- **Sparsity**: Dimensions depends on vocabulary length, sparse vectors, inefficient for computation.
- Lack of Context: Do not capture the context or semantics of words, limiting the representation's expressiveness.
- Synonymy and Polysemy: Struggle with words that have multiple meanings or similar meanings, leading to ambiguity in representation.

Introduction & Motivations: Capturing Semantic Similarity

Semantic Similarity - A Human Concept:

- Traditional models often fail to capture the inherent semantic similarity between concepts that are intuitively understood by humans.
- Examples of semantically similar concepts:
 - 'stag' and 'deer'
 - 'osteopath' and 'physiotherapist'
 - 'prince' and 'king'
- This gap highlights the need for models that can effectively *transfer* human-like understanding of language and semantics.





Deer Credits: DuckDuckGo Research

Semantic Similarity with TF-IDF

TF-IDF Representation:

- Vocabulary: {'cat', 'dog', 'apple', 'car', 'tree', 'robot'}
- *w_i* : *x_j* = 1 if *j* = *i* else 0

Limitation - Words Distance

- In this high-dimensional space, distance between words means nothing.
- Fails to capture any semantic or contextual relationships between words.



High-dimensional space of Tf-IDF Vectorizer

Word Embeddings

Concept of Word Embeddings

Ideal representations: Vector representations of words in a vector space where semantically similar words are mapped to nearby points.

Concept of Word Embeddings

Ideal representations: Vector representations of words in a vector space where semantically similar words are mapped to nearby points. **Key Properties:**

 Semantic Relationships: Words with similar meanings are located in close proximity. The closer the words, the more similar.



Ideal Vector Space Credits: Wikipedia

Concept of Word Embeddings

Ideal representations: Vector representations of words in a vector space where semantically similar words are mapped to nearby points. **Key Properties:**

- Semantic Relationships: Words with similar meanings are located in close proximity. The closer the words, the more similar.
- Syntactic Relationships: Capturing patterns in word use based on the context, and certain algebraic "operations" to produce meaningful relationships (e.g., "king" - "man" + "woman" "queen").





Evaluations of the Embeddings

Evaluating Word Representation Word representation effectiveness can be assessed through intrinsic and extrinsic evaluations. **Intrinsic Evaluation:** Measures how well the embedding captures linguistic properties:

 Word Similarities: Evaluating the closeness of words in the embedding space



Figure 2: Two-dimensional PCA projection of the 1000-dimensional Skip-gram vectors of countries and their copilal cities. The figure libratures ability of the model to automatically organize concepts and learn implicitly the relationships between them, as during the training we did not provide any supervised information about what a capital city means.

Evaluating Word Representation Word representation effectiveness can be assessed through intrinsic and extrinsic evaluations. **Intrinsic Evaluation:** Measures how well the embedding captures linguistic properties:

- Word Similarities: Evaluating the closeness of words in the embedding space
- Word Analogies: Testing the embedding's ability to deduce relationships.



Figure 2: Two-dimensional PCA projection of the 1000-dimensional Skip-gram vectors of countries and their capital cities: The figure illustrates ability of the model to automatically organize concepts and learn implicitly the relationships between them, as during the training we did not provide any supervised information about what a capital city means.

Evaluating Word Representation Word representation effectiveness can be assessed through intrinsic and extrinsic evaluations. **Intrinsic Evaluation:** Measures how well the embedding captures linguistic properties:

- Word Similarities: Evaluating the closeness of words in the embedding space
- Word Analogies: Testing the embedding's ability to deduce relationships.
- Synonym Detection: Identifying words with similar meanings.



Figure 2: Two-dimensional PCA projection of the 1000-dimensional Skip-gram vectors of countries and their capital cities: The figure illustrates ability of the model to automatically organize concepts and learn implicitly the relationships between them, as during the training we did not provide any supervised information about what a capital city means.

Evaluating Word Representation Word representation effectiveness can be assessed through intrinsic and extrinsic evaluations. **Intrinsic Evaluation:** Measures how well the embedding captures linguistic properties:

- Word Similarities: Evaluating the closeness of words in the embedding space
- Word Analogies: Testing the embedding's ability to deduce relationships.
- Synonym Detection: Identifying words with similar meanings.
- Word Clustering: Grouping semantically similar words.



Figure 2: Two-dimensional PCA projection of the 1000-dimensional Skip-gram vectors of countries and their capital cities: The figure illustrates ability of the model to automatically cognize zocorcepts and learn implicit the relationships between them, as during the training we did not provide any supervised information about what a capital city means.

Intrinsic Evaluation Example: Word Similarities

Assessing Word Similarities:

- Evaluates how well the embedding captures semantic similarities between words.
- Measures the cosine similarity between vectors representing different words.

Example:

- Comparing words such as 'king' and 'queen' versus 'king' and 'apple'.
- Expect 'king' and 'queen' to have a higher similarity score than 'king' and 'apple'.

Importance: Understanding word similarities allows for a nuanced understanding of the embedding space, reflecting how well the model captures semantic relationships.



Word similarities

Named Entity Recognition (NER):

Identifying named entities in text.

1		1	TS Teach - Calogle	n	Ľ	88.7	30.8	87.5	91 989 2	K7.5/10.1	76.8/98.4
1		2	ALBERT Team Cougle Langua	pAURIT (Insentia)	ß	35.4	60.3	81.5	93.491.2	\$2,5/92.0	34,299.5
-	F.	3	1 3	AUCE v2 large ensemble (Albaba DAMO NUP)	Ľ	85.0	69.2	87.1	93.691.5	\$2,792,3	34,490.1
1		4	Microsoft 2005 A/ & UMD	FixeLD-RoBERTs (ensemble)	Ľ,	05.0	60.8	96.8	93.190.8	82.492.2	74,898.0
1		5	Facebook Al	Astern	Ľ	86.5	62.8	96.2	92.209.8	R2 2/54.9	34398.2
1		6	XLMH Than	X3.76H Carge (investible)	ß	38.4	87.8	86.8	93.0190.7	\$1.876.1	24.2/98.3
+	ŧ.	τ	Mexical 2005 ALLMSR N	NT GAN ersenkle	Ø	87.6	65.4	96.5	92 1190.3	\$1,090.7	73.758.0
		ð	GLUE Human Baselines	GLUE Human Baselines	ß	87.1	65.4	\$7.8	86.5100.8	\$2,792.6	59.558.4
		9	Station Hazy Research	Section Mellat.	Ľ	03.2	63.8	96.2	91500.5	90.1/09.7	73.508.9
1		20	XLM Systems	XLM (English-sely)	Ľ	83.1	62.8	85.5	96797.5	80.859.2	72,269.0

- Named Entity Recognition (NER): Identifying named entities in text.
- **Text Classification**: Categorizing text into predefined classes.

	1	TS Teach - Gaogle	n	Ľ	88.7	30.8	87.5	91 989 2	KZ 5/102 1	21.8/98.4
	z	ALBERT from Coople Langua	(element)	ß	35.4	69.3	87.5	93.491.2	\$2,5/92,0	34,299.5
÷	3	1 /8	AUCE v2 large ensemble (Albaba DAMO HLP)	ß	85.0	69.2	87.1	93.691.5	\$2,792,3	74.490.1
		Microsoft DB55 A/ & UMD	FixeLD-RobERTs (insemble)	Ľ	06.0	60.8	96.8	93.1.90.8	82.492.2	74,898.0
	5	Facebook A/	Roberts	Ľ	86.5	62.8	96.2	92.209.8	82.2554.9	34.398.2
	6	XLMM Train	XL/en Large (maerilike)	Ľ	38.4	17.8	86.8	93.0190.7	81.8/06.1	34,2/98.3
÷	τ	Mexical 2005 ALL MSR.N	NT 04A ensemble	ß	87.6	65.4	96.5	92 790.3	\$1,1/90.7	73.758.0
	ō	GLUE Human Baselines	GLUE Haman Baselines	ß	87.1	65.4	\$7.8	86.5100.8	\$2,792.6	50.5/58.4
	9	Stanford Hazy Research	Sected Mellat.	Ľ	03.2	63.8	96.2	91500.5	901/09.7	73.308.9
	20	XLM Systems	XLM (English-sely)	ß	83.1	62.8	85.5	96797.5	80.859.2	72,269.0

- Named Entity Recognition (NER): Identifying named entities in text.
- **Text Classification**: Categorizing text into predefined classes.
- Part-of-Speech (POS) Tagging: Labeling words with their corresponding part of speech.

	1	TS Teach - Gaogle	n	Ľ	88.7	30.8	87.5	91 989 2	K7.5/10.1	21.8/98.4
	2	ALBERT from Coople Langua	(element)	ß	35.4	69.3	87.5	93.491.2	\$2,5/52,0	34,299.5
÷	3	1 /8	AUCE v2 large ensemble (Albaba DAMO HLP)	ß	85.0	69.2	\$7.1	93.691.5	\$2,792,3	74.490.1
	.4	Microsoft DB55 A/ & UMD	FixeLD-RobERTs (insemble)	Ľ	05.0	60.8	96.8	93.1.90.8	82.492.2	74,898.0
	5	Facebook A/	Roberts	Ľ	18.5	62.8	96.2	92.209.8	82.2554.9	34.398.2
	6	XLMM Train	XL/en Large (maerilike)	Ľ	38.4	17.8	86.8	93.0190.7	\$1.876.1	34,2/98.3
÷		Mexical 2005 ALL MSR.N	NT 04A ensemble	ß	87.6	65.4	96.5	92 790.3	\$1,090.7	73.758.0
		GLUE Human Baselines	GLUE Haman Baselines	ß	87.1	65.4	\$7.8	86.5100.8	\$2,792.6	50.5/58.4
	9	Stanford Hazy Research	Sected Mellat.	Ľ	03.2	63.8	96.2	91500.5	90.1/09.7	73.308.9
	22	XLM Systems	XLM (English-sely)	ß	83.1	62.8	85.5	96797.5	80.859.2	72,269.0

- Named Entity Recognition (NER): Identifying named entities in text.
- **Text Classification**: Categorizing text into predefined classes.
- Part-of-Speech (POS) Tagging: Labeling words with their corresponding part of speech.
- ...and other NLP tasks.

I											
I		1	TS Team - Gaogle	n	Ľ	88.7	30.8	87.5	91,989,2	82.5/92.1	24.8/98.4
		2	ALBERT Team Cougle Language	pAURIT (Crossic)	ß	88.4	60.3	87.5	93.491.2	\$2,5/02.0	34,2/98.1
ł	٠	3	1 /8	AUCE v2 large ensemble (Albaba DAMO HLP)	ß	85.0	69.2	\$7.1	93.691.5	\$2,792,3	34,490.1
		.4	Microsoft D005 A/ & UMD	FixeLD-RobERTa (intertble)	ß	05.0	60.8	96.8	93.190.8	\$2,492.2	74,898.5
		5	Facebook Al	Astern	Ľ	86.5	62.8	96.2	92.209.8	K2.2154.9	343983
		6	XLMM Train	XL/en Large (maerilike)	Ľ	38.4	17.8	86.8	93.0190.7	81.8/06.1	34,2/98.3
ł	÷		Mexical 2005 ALL MSR.N	NT 04A ensemble	ß	87.6	65.4	96.5	92 790.3	\$1,1/90.7	73.758.0
		ð	GLUE Human Baselines	GLUE Haman Baselines	ß	87.1	65.4	\$7.8	86.5100.8	\$2,792.6	59.558.4
		9	Stanford Hazy Research	Sected Mellat.	Ľ	03.2	63.8	96.2	91500.5	901/09.7	73.508.5
		20	XLM Systems	XLM (English-sely)	ß	83.1	62.8	85.5	96797.5	80.859.2	72,259.6

Extrinsic Evaluation Example: Text Classification

Text Classification with Embeddings:

- Evaluates the effectiveness of word embeddings in text classification.
- Measures the improvement in classification accuracy when using embeddings.

Example:

- Using word embeddings as features for a sentiment analysis model.
- Comparing model performance with and without the use of embeddings.

Importance: Demonstrates the practical utility of word embeddings in real-world applications, highlighting their contribution to model performance.

Architecture	Accuracy [%]					
4-gram [32]	39					
Average LSA similarity [32]	49					
Log-bilinear model [24]	54.8					
RNNLMs [19]	55.4					
Skip-gram	48.0					
Skip-gram + RNNLMs	58.9					

Microsfot Sentence

Completion Challenge

Credit: Mikolov et al. 2013

One proposal

Ontologies

Structured framework to define and categorize entities within a domain.

- Level: Distinguishing between instances.
- Class: Grouping entities into collections or concepts.
- Attribute: Describing entities through characteristics.
- ...and more domain-specific categorizations.

Domain-Specific Ontologies:

- Medicine: ICD-9 or ICD-10 for diseases.
- Computer Science: Structured codes.

Limitations: Ontologies are not scalable: establishing comprehensive links between all entities can be labor-intensive and complex.



Word Representation through Ontology Credits: Jay Alammar

Continuous Vector representation

Dense embeddings provide compact, rich representations of words, capturing semantic and syntactic nuances effectively.

- Key Benefits:
 - Semantic Richness: Encapsulate meanings.



Dense embeddings provide compact, rich representations of words, capturing semantic and syntactic nuances effectively.

Key Benefits:

- Semantic Richness: Encapsulate meanings.
- Reduced Dimensionality: Lower-dimensional space.



Dense embeddings provide compact, rich representations of words, capturing semantic and syntactic nuances effectively.

Key Benefits:

- Semantic Richness: Encapsulate meanings.
- Reduced Dimensionality: Lower-dimensional space.
- Mathematical Operability: Enable operations like analogy solving.



Dense embeddings provide compact, rich representations of words, capturing semantic and syntactic nuances effectively.

Key Benefits:

- Semantic Richness: Encapsulate meanings.
- Reduced Dimensionality: Lower-dimensional space.
- Mathematical Operability: Enable operations like analogy solving.
- Enhanced Performance: Boost NLP tasks with contextually-aware representations.



Dense embeddings provide compact, rich representations of words, capturing semantic and syntactic nuances effectively.

Key Benefits:

- Semantic Richness: Encapsulate meanings.
- Reduced Dimensionality: Lower-dimensional space.
- Mathematical Operability: Enable operations like analogy solving.
- Enhanced Performance: Boost NLP tasks with contextually-aware representations.



Dense embeddings provide compact, rich representations of words, capturing semantic and syntactic nuances effectively.

Key Benefits:

- Semantic Richness: Encapsulate meanings.
- Reduced Dimensionality: Lower-dimensional space.
- Mathematical Operability: Enable operations like analogy solving.
- Enhanced Performance: Boost NLP tasks with contextually-aware representations.

Conclusion: Dense embeddings represent a leap in language modeling, facilitating advanced architectures and deeper language understanding.



Word2Vec - Mikolov et al. (2013): Word2Vec offers two architecture choices for generating dense word embeddings, inspired by language modeling: Continuous Bag-Of-Words (CBOW):

Represents words through n-gram context.

Continuous Skip-gram:



Word2Vec - Mikolov et al. (2013): Word2Vec offers two architecture choices for generating dense word embeddings, inspired by language modeling: Continuous Bag-Of-Words (CBOW):

- Represents words through n-gram context.
- Uses projection matrices (embedding layers) to capture word features.

Continuous Skip-gram:



Word2Vec - Mikolov et al. (2013): Word2Vec offers two architecture choices for generating dense word embeddings, inspired by language modeling: Continuous Bag-Of-Words (CBOW):

- Represents words through n-gram context.
- Uses projection matrices (embedding layers) to capture word features.
- Aims to predict a target word based on surrounding context words.

Continuous Skip-gram:



Word2Vec - Mikolov et al. (2013): Word2Vec offers two architecture choices for generating dense word embeddings, inspired by language modeling: Continuous Bag-Of-Words (CBOW):

- Represents words through n-gram context.
- Uses projection matrices (embedding layers) to capture word features.
- Aims to predict a target word based on surrounding context words.

Continuous Skip-gram:

Also uses n-gram representation and projection matrices.



Word2Vec - Mikolov et al. (2013): Word2Vec offers two architecture choices for generating dense word embeddings, inspired by language modeling: Continuous Bag-Of-Words (CBOW):

- Represents words through n-gram context.
- Uses projection matrices (embedding layers) to capture word features.
- Aims to predict a target word based on surrounding context words.

Continuous Skip-gram:

- Also uses n-gram representation and projection matrices.
- In contrast to CBOW, predicts surrounding context words from a target word, offering quality embeddings for even infrequent words.


Training the Skip-Gram Model: Data Collection

Objective of Skip-Gram: The Skip-Gram model aims to predict context words given a target word, strengthening the word associations within a specified window size.

Objective of Skip-Gram: The Skip-Gram model aims to predict context words given a target word, strengthening the word associations within a specified window size.

Training Data Preparation:

- Features (Input): Target words.
- Labels (Output): Context words within a defined window size around the target word.

Objective of Skip-Gram: The Skip-Gram model aims to predict context words given a target word, strengthening the word associations within a specified window size.

Training Data Preparation:

- Features (Input): Target words.
- Labels (Output): Context words within a defined window size around the target word.

We don't need labelled data: we can collect texts from Wikipedia, books, internet... and create the training set from it !

Example: Given the sentence "The quick brown fox jumps over", with a window size of 1, the training pairs are:

- Input: "quick", Output: ["The", "brown"]
- Input: "brown", Output: ["quick", "fox"]
- ... and so on.



Training pairs generation in Skip-Gram Model

Model Architecture:

- Input Layer: One-hot encoded vectors of the target words.
- Embedding layer: A fully connected layer of dimension
 |Vocabulary| × Embedding size n
- **Hidden Layer**: A fully connected layer without activation, to project to the output vocabulary.
- Output Layer: Predicts the probability distribution (softmax) of context words for the given input.

Training Skip-Gram Model: Model Architecture (2/2)

Skip-Gram Model Architecture



Skip-Gram Model Architecture; predict context from one word

Main Limitation: each time we're computing softmax for V words (could be 10e6!) the process is computationally expensive..

Challenge in Training: The standard Skip-Gram model with softmax can be computationally expensive due to the large vocabulary size.

Challenge in Training: The standard Skip-Gram model with softmax can be computationally expensive due to the large vocabulary size. **Solution - Negative Sampling:**

• **Concept**: Instead of predicting the probability for all words in the vocabulary, negative sampling trains the model to distinguish a target word from a small set of random 'noise words'.

Challenge in Training: The standard Skip-Gram model with softmax can be computationally expensive due to the large vocabulary size. **Solution - Negative Sampling:**

- **Concept**: Instead of predicting the probability for all words in the vocabulary, negative sampling trains the model to distinguish a target word from a small set of random 'noise words'.
- Benefits: Simplifies the computation and accelerates the training process, especially for large corpora.

Challenge in Training: The standard Skip-Gram model with softmax can be computationally expensive due to the large vocabulary size. **Solution - Negative Sampling:**

- **Concept**: Instead of predicting the probability for all words in the vocabulary, negative sampling trains the model to distinguish a target word from a small set of random 'noise words'.
- Benefits: Simplifies the computation and accelerates the training process, especially for large corpora.

Challenge in Training: The standard Skip-Gram model with softmax can be computationally expensive due to the large vocabulary size. **Solution - Negative Sampling:**

- Concept: Instead of predicting the probability for all words in the vocabulary, negative sampling trains the model to distinguish a target word from a small set of random 'noise words'.
- Benefits: Simplifies the computation and accelerates the training process, especially for large corpora.

Implementation: In each training step:

Select a small number of negative samples (words not in the context).

Challenge in Training: The standard Skip-Gram model with softmax can be computationally expensive due to the large vocabulary size. **Solution - Negative Sampling:**

- **Concept**: Instead of predicting the probability for all words in the vocabulary, negative sampling trains the model to distinguish a target word from a small set of random 'noise words'.
- Benefits: Simplifies the computation and accelerates the training process, especially for large corpora.

Implementation: In each training step:

- Select a small number of negative samples (words not in the context).
- Update weights based on the target word and the sampled negative words.

Example: Given the sentence "The quick brown fox jumps over", with a window size of 1, the training pairs are:

- Input: ["quick", "brown"] Output: 1
- Input: ["quick", "yellow"] Output 0
- Input: ["brown", "fox"] Output: 1
- Input: ["brown", "fly"] Output: 0
- ... and so on.

Skip-Gram Input-Output Pairs				
quick	The	1		
quick	brown	1		
quick	bow	0		
quick	yellow	0		
brown	quick	1		
brown	fox	1		
brown	wolf	0		
brown	maths	0		
fox	brown	1		
fox	jumps	1		
fox	bark	0		
fox	fly	0		
jumps	fox	1		
jumps	over	1		
jumps	dig	0		
jumps	special	0		
Input	Output	Value		

Training pairs generation in Skip-Gram Model with Negative Sampling

Model Architecture:

- Input Layer: One-hot encoded vectors of the target words. One-hot Encoded vectors of the context words.
- Embedding layer: A fully connected layer of dimension
 |Vocabulary| × Embedding size n
- Context Layer: A fully connected layer of dimension |Vocabulary| × Embedding size n
- Output Layer: Predicts the probability distribution (softmax) of second word to be the context of first word.

Training Skip-Gram Model with NS: Model Architecture (2/2)

Skip-Gram Model Architecture





Skip-Gram Model Architecture; predict context from one word

Main Difference: each time we're computing softmax for one true target and some noise instead of the vocabulary *V* words !

- Cosine Similarity: Measures the cosine of the angle between two word vectors, indicating how similar they are in the embedding space.
- Word pair relationship analysis has practical uses in fields like semantic search, automated text analysis, and even in creative domains like generating novel content based on identified patterns.
- A deeper understanding of word relationships can enhance language models, making them more robust and contextually aware.

Table 8: Examples of the word pair relationships, using the best word vectors from Table 4 (Skipgram model trained on 783M words with 300 dimensionality).

Relationship	Example 1	Example 2	Example 3
Terminoniship	Estample 1	Entample B	Estumple 5
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

Word Pairs Relationships Credit: Mikolov et al. (2013)

GloVe - Global Vectors for Word Representation:

- **Co-occurrence Matrix**: GloVe is built on word-word co-occurrence statistics from a corpus, capturing global statistical information.
- **Applications**: Widely used in applications requiring an understanding of word similarity and analogy based on global corpus statistics.

FastText - Advanced Word Representation:

- **Subword Information**: Extends the Word2Vec model by representing each word as a bag of character n-grams, capturing morphological information.
- Handling Rare Words: Particularly effective in understanding and representing rare words or misspellings.
- **Applications**: Useful in tasks where morphological information is crucial, like language modeling and text classification in morphologically rich languages.

 Word Sense Ambiguity: only one vector per word, ignoring the polysemy where words have multiple meanings based on context.

- Word Sense Ambiguity: only one vector per word, ignoring the polysemy where words have multiple meanings based on context.
- Context Ignorance: Unable to capture the meaning of a word in different contexts. The same word in different sentences will have the same representation.

- Word Sense Ambiguity: only one vector per word, ignoring the polysemy where words have multiple meanings based on context.
- Context Ignorance: Unable to capture the meaning of a word in different contexts. The same word in different sentences will have the same representation.
- Out-of-Vocabulary (OOV) Words: Challenges in handling words not present in the training corpus. FastText addresses this partially with subword information.

- Word Sense Ambiguity: only one vector per word, ignoring the polysemy where words have multiple meanings based on context.
- Context Ignorance: Unable to capture the meaning of a word in different contexts. The same word in different sentences will have the same representation.
- Out-of-Vocabulary (OOV) Words: Challenges in handling words not present in the training corpus. FastText addresses this partially with subword information.
- Fixed Representations: Fixed after training, not allowing the model to adapt to evolving language use or domain-specific jargon.

- Word Sense Ambiguity: only one vector per word, ignoring the polysemy where words have multiple meanings based on context.
- Context Ignorance: Unable to capture the meaning of a word in different contexts. The same word in different sentences will have the same representation.
- Out-of-Vocabulary (OOV) Words: Challenges in handling words not present in the training corpus. FastText addresses this partially with subword information.
- Fixed Representations: Fixed after training, not allowing the model to adapt to evolving language use or domain-specific jargon.
- Resource Intensive: Requires substantial computational resources and time to train on large corpora, making it less feasible for resource-constrained scenarios.

Contextual Embeddings

Overcoming the Limitations: Contextual embeddings represent the next evolution in word representations, addressing the inherent constraints of static word embeddings.

What are Contextual Embeddings?

- Dynamic Word Representations: Unlike static embeddings, contextual embeddings provide representations that change based on the word's context.
- Deep Contextualization: These models consider the entire sentence or even larger contexts to understand the meaning of each word.

Pioneering Models:

- ELMo (Embeddings from Language Models), Peters et al. (2018): Utilizes bidirectional LSTM trained on a specific task to generate embeddings.
- BERT (Bidirectional Encoder Representations from Transformers), Delvin et al. (2019): Transforms the landscape with a transformer-based model, pre-trained on vast amounts of text and fine-tuned for specific tasks.
- GPT-2 (Generative Pre-trained Transformer), Radford et al. (2019): Emphasizes generative capabilities and large-scale pre-training for versatile language understanding.

Advantages Over Static Embeddings:

- Captures polysemy by providing context-specific word meanings.
- Adapts to different domains and evolving language use without retraining from scratch.
- Enhances performance across a wide array of NLP tasks.

Conclusion: Contextual embeddings mark a significant milestone in NLP, offering nuanced and adaptable understanding of language, far surpassing the capabilities of static embeddings.

What is ELMo?

- ELMo, Peters et al. (2018) was developed by Allen Institute for Al.
- It stands for Embeddings from Language Models.

Key Features of ELMo:

- Deep Contextualization: ELMo considers the entire context of a word by using the internal states of a bidirectional LSTM.
- **Dynamic Word Representation**: Each word's representation is a function of the entire sentence.

Architecture Overview: ELMo combines the representations of a two-layer pretrained bidirectional LSTM internal states into downstream tasks to outperfom current models.

Pre-training Procedure:

- Bidirectional Language Modeling:
 - Forward LSTM: Predicts the next word from the past context.
 - *Backward LSTM*: Predicts the previous word from the future context.
- Character Embeddings: Captures word morphology and manages out-of-vocabulary words.
- **Training on Large Corpus**: Enhances understanding of language structure by predicting surrounding words.

Objective: ELMo's pre-training on a large corpus with bidirectional context and character embeddings lays a robust foundation for nuanced language understanding.

Forward LM LSTM training

Forward LM with LSTMs



Forward LSTM

Architecture:

- One embedding layer.
- Two hidden layers.
- One softmax layer to predict next token.

Backward LM LSTM training



Backward LM with LSTMs

Backward LSTM

Architecture:

- One embedding layer.
- Two hidden layers.
- One softmax layer to predict previous token.

Final Embeddings

Final Representation



Final Embeddings

Architecture:

- We concatenate each layer together.
- We weight each of the layer to get the final embedding.
- The weights depend of the downstream task.

 Feature Extraction: Use the output of the pre-trained LSTM layers as features in a new model tailored to a specific task, such as sentiment analysis or named entity recognition.

- Feature Extraction: Use the output of the pre-trained LSTM layers as features in a new model tailored to a specific task, such as sentiment analysis or named entity recognition.
- **Fine-tuning**: Adjust the pre-trained LSTM weights slightly during the training of the downstream task to better adapt the model to the specific requirements of the task.

- Feature Extraction: Use the output of the pre-trained LSTM layers as features in a new model tailored to a specific task, such as sentiment analysis or named entity recognition.
- **Fine-tuning**: Adjust the pre-trained LSTM weights slightly during the training of the downstream task to better adapt the model to the specific requirements of the task.

- Feature Extraction: Use the output of the pre-trained LSTM layers as features in a new model tailored to a specific task, such as sentiment analysis or named entity recognition.
- **Fine-tuning**: Adjust the pre-trained LSTM weights slightly during the training of the downstream task to better adapt the model to the specific requirements of the task.

Downstream Task Architecture:

• Start with the pre-trained ELMo embeddings as the input layer.

- Feature Extraction: Use the output of the pre-trained LSTM layers as features in a new model tailored to a specific task, such as sentiment analysis or named entity recognition.
- **Fine-tuning**: Adjust the pre-trained LSTM weights slightly during the training of the downstream task to better adapt the model to the specific requirements of the task.

Downstream Task Architecture:

- Start with the pre-trained ELMo embeddings as the input layer.
- Add task-specific layers on top of the ELMo layers (e.g., additional LSTM layers, dense layers...).
Integration into Downstream Tasks:

- Feature Extraction: Use the output of the pre-trained LSTM layers as features in a new model tailored to a specific task, such as sentiment analysis or named entity recognition.
- **Fine-tuning**: Adjust the pre-trained LSTM weights slightly during the training of the downstream task to better adapt the model to the specific requirements of the task.

Downstream Task Architecture:

- Start with the pre-trained ELMo embeddings as the input layer.
- Add task-specific layers on top of the ELMo layers (e.g., additional LSTM layers, dense layers...).
- The final output layer is designed according to the downstream task (e.g., softmax for classification).

ELMo embeddings can enhance the sentiment analysis models by providing deep contextualized word representations.

Incorporation into Model (example):

Input Layer: Start with ELMo embeddings for each token in the input text.

ELMo embeddings can enhance the sentiment analysis models by providing deep contextualized word representations.

- Input Layer: Start with ELMo embeddings for each token in the input text.
- Additional Layers: Add a bi-directional LSTM layer to further capture context not encapsulated by ELMo.

ELMo embeddings can enhance the sentiment analysis models by providing deep contextualized word representations.

- Input Layer: Start with ELMo embeddings for each token in the input text.
- Additional Layers: Add a bi-directional LSTM layer to further capture context not encapsulated by ELMo.
- **Output Layer**: A dense layer with softmax activation to classify the sentiment as positive or negative.

ELMo embeddings can enhance the sentiment analysis models by providing deep contextualized word representations.

- Input Layer: Start with ELMo embeddings for each token in the input text.
- Additional Layers: Add a bi-directional LSTM layer to further capture context not encapsulated by ELMo.
- Output Layer: A dense layer with softmax activation to classify the sentiment as positive or negative.
- Train your model: The whole model with ELMo frozen or not !

ELMo embeddings can enhance the sentiment analysis models by providing deep contextualized word representations.

- Input Layer: Start with ELMo embeddings for each token in the input text.
- Additional Layers: Add a bi-directional LSTM layer to further capture context not encapsulated by ELMo.
- Output Layer: A dense layer with softmax activation to classify the sentiment as positive or negative.
- Train your model: The whole model with ELMo frozen or not !

ELMo embeddings can enhance the sentiment analysis models by providing deep contextualized word representations.

Incorporation into Model (example):

- Input Layer: Start with ELMo embeddings for each token in the input text.
- Additional Layers: Add a bi-directional LSTM layer to further capture context not encapsulated by ELMo.
- **Output Layer**: A dense layer with softmax activation to classify the sentiment as positive or negative.
- Train your model: The whole model with ELMo frozen or not !

Advantage: The use of ELMo *transfers* the nuances of language, leading to a more accurate sentiment classification compared to models without contextualized embeddings: **3.3% in absolute improvement**.

ELMo's Impact on NLP Tasks

ELMo's introduction marked a new era in NLP by setting state-of-the-art (SOA) benchmarks across multiple tasks simultaneously.

Remarkable Achievements:

- SOA in Six Benchmarks: ELMo established new records in a range of NLP tasks.
- Double-Digit Improvements: Notably increased performance by over 10% in four of those tasks.

Transfer Learning with ELMo:

- Knowledge Acquisition: Gained from pre-training on extensive datasets.
- Knowledge Transfer: Applied to enhance task-specific models, demonstrating the power of transfer learning in NLP.

	TASK	PREVIOUS SOTA		OUR BASELINE	ELMO + BASELINE	(ABSOLUTE/ RELATIVE)
Q&A	SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
Textual entailment	SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17	0.7/5.8%
Semantic role labelling	SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coreference resolution	Coref	Lee et al. (2017)	67.2	67.2	70.4	3.2/9.8%
amed entity recognition	NER	Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10	2.06/21%
Sentiment analysis	SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5	3.3/6.8%

Table 1: Test set comparison of ELMo enhanced neural models with state-of-the-art single model headings are site benchmark. NLP tasks. The performance metric varies across tasks – accuracy for SNLI and STS-5; Fi 5 SQAAD, SRL and NER; serenge F; for Coref. Date to the small test size for NER and SST-5, we report the meta and standard deviation across live runs with different random seeds. The "increase" oclarm lists both the absdut and relative improvements over our baseline.

Visualization of ELMo generating embeddings Credit: Pieters Benefits of Using ELMo's LSTM:

- Contextualized Understanding: Models become more aware of the context within sentences, leading to better understanding and predictions.
- **Transfer Learning**: Leveraging pre-trained models significantly reduces the need for large labeled datasets for the target task.
- State-of-the-Art Results: Many tasks see improved performance metrics when integrating pre-trained LSTMs from ELMo.

Related Work in Contextual Embeddings

Several works have built upon and extended the concepts introduced by ELMo, pushing the boundaries in various NLP tasks.

- ULM-Fit (Howard and Ruder, 2018): A model that employs transfer learning specifically for text classification, setting new benchmarks on six different tasks.
- Semi-supervised Sequence Tagging (Peters et al., 2017): Improved named entity recognition (NER) performance on multiple datasets by using bidirectional language models.
- Learned in Translation: Contextualized Word Vectors (McCann et al., 2017): Achieved state-of-the-art results in translation tasks by leveraging contextualized word vectors derived from an LSTM model with an attention mechanism.

These works collectively demonstrate the growing impact of transfer learning and contextual embeddings across a range of NLP applications. 37

Challenges in Cutting-Edge NLP: Models like ELMo and ULM-Fit push the boundaries but encounter challenges:

- Data and Resources: Dependence on extensive datasets and substantial computational power limits accessibility and applicability, especially for under-resourced languages and domains.
- Model Complexity: The intricate nature of these models can obscure interpretability, leading to the 'black-box' issue and challenges in model trust and fine-tuning.
- Real-World Application: Mastery over benchmarks doesn't always translate to real-world scenarios, where data can be noisy and unpredictable.

Conclusion: The future of NLP lies in overcoming these hurdles, striving for models that balance performance with efficiency, interpretability, and broader applicability. 38

Open Discussion

- Feel free to ask questions or share your thoughts about today's topics.
- Any insights, experiences, or perspectives you'd like to discuss are welcome.

Summary of Key Takeaways

- We explored static embeddings like Word2Vec, GloVe, and FastText, and their role in establishing the foundation for current NLP advancements.
- Static embeddings, while transformative, have limitations in handling polysemy, dynamic context, and out-of-vocabulary words.
- Advanced NLP models like ELMo and ULM-Fit build upon static embeddings, offering context-aware representations that better capture the nuances of language.
- These advanced models set new benchmarks but face challenges regarding data dependency, computational demands, and generalization.
- Addressing the limitations of both static and advanced models is crucial for the development of more efficient, interpretable, and generalizable NLP solutions.