# Contents

# WEEK 6: FROM AUTOENCODERS TO EMBEDDINGS

## Understanding Word Embeddings

Word embeddings are a fundamental concept in natural language processing (NLP) that allows words to be represented as dense vectors in a high-dimensional space. These vector representations capture semantic and syntactic relationships between words, enabling machines to better understand and process human language. One of the most popular techniques for generating word embeddings is Word2Vec, introduced by Mikolov et al. in 2013. Word2Vec is an unsupervised learning algorithm that learns word embeddings from large corpora of text data.

Word2Vec offers two distinct architectures for learning word embeddings: the Continuous Bag-of-Words (CBOW) model and the Skip-gram model. The CBOW model predicts the target word given its surrounding context words, while the Skip-gram model predicts the context words given the target word. In the CBOW model, the objective is to maximize the probability of the target word given the average of the embeddings of its context words. Mathematically, this can be expressed as:

$$p(w_t|w_{t-k},...,w_{t-1},w_{t+1},...,w_{t+k}) = \frac{exp(v_{w_t}^T \cdot \frac{1}{2k}\sum_{i=1}^{2k} v_{w_{t+i}})}{\sum_{w \in V} exp(v_w^T \cdot \frac{1}{2k}\sum_{i=1}^{2k} v_{w_{t+i}})}$$

where $w_t$ is the target word, $w_{t-k},...,w_{t-1},w_{t+1},...,w_{t+k}$ are the context words, $v_w$ is the embedding vector for word $w$, and $V$ is the vocabulary.

In contrast, the Skip-gram model aims to maximize the probability of the context words given the target word. The objective function for the Skip-gram model is:

$$p(w_{t-k},...,w_{t-1},w_{t+1},...,w_{t+k}|w_t) = \prod_{i=1}^{2k} p(w_{t+i}|w_t)$$

where $p(w_{t+i}|w_t)$ is computed using the softmax function:

$$p(w_o|w_i) = \frac{exp(v_{w_o}^T \cdot v_{w_i})}{\sum_{w \in V} exp(v_w^T \cdot v_{w_i})}$$

The context window is a crucial hyperparameter in Word2Vec that determines the number of words considered as context for each target word. A larger context window captures more global information, while a smaller context window focuses on local relationships. The choice of context window size depends on the specific task and the desired balance between capturing local and global word relationships.

To improve the efficiency of training Word2Vec models, negative sampling is often employed. Negative sampling is a simplified version of noise contrastive estimation (NCE) that reduces the computational

complexity of the softmax function. Instead of considering all words in the vocabulary as potential context words, negative sampling selects a small number of "negative" samples (i.e., words that are unlikely to appear in the context of the target word) and trains the model to distinguish between the true context words and the negative samples. This approach significantly reduces the number of computations required during training, making it feasible to learn word embeddings from large datasets.

## Properties of Word Embeddings

Word embeddings are dense vector representations of words that capture semantic and syntactic information from large corpora of text. These embeddings are typically learned using neural network architectures, such as the Skip-gram or Continuous Bag-of-Words (CBOW) models, which aim to predict a target word given its context or vice versa. The resulting vectors exhibit several interesting properties that make them useful for various natural language processing tasks.

One of the most notable properties of word embeddings is their ability to capture semantic relationships between words. Words with similar meanings tend to have vectors that are close to each other in the embedding space. This proximity can be measured using cosine similarity, which calculates the cosine of the angle between two vectors. For example, the vectors for words like "cat" and "dog" would be closer to each other than the vectors for "cat" and "airplane," reflecting their semantic similarity. This property allows word embeddings to be used for tasks such as word similarity, word clustering, and semantic search.

Another fascinating property of word embeddings is the ability to perform arithmetic operations on word vectors to uncover relationships and analogies. By adding and subtracting word vectors, we can explore the relationships between words and concepts. For instance, the famous example "king - man + woman = queen" demonstrates that the difference between the vectors for "king" and "man" represents the concept of gender, and when added to the vector for "woman," it results in the vector for "queen." This arithmetic property can be generalized to other relationships, such as "capital - country + country = capital," where the difference between a country's vector and its capital's vector represents the concept of a capital city. These arithmetic operations provide a powerful tool for understanding and manipulating semantic relationships within the embedding space.

Word embeddings can also be used to solve analogy tasks, which involve finding a word that shares the same relationship with a given word as another pair of words. The classic example is "king : man :: queen : woman," where the goal is to find the word that relates to "man" in the same way that "queen" relates to "woman." To solve this analogy using word embeddings, we can leverage the arithmetic property mentioned earlier. By calculating "king - man + woman," we obtain a vector that is close to the vector for "queen," effectively solving the analogy. This approach can be applied to various types of analogies, such as "Paris : France :: Tokyo : Japan," demonstrating the versatility of word embeddings in capturing different types of relationships.

The ability of word embeddings to capture semantic relationships, perform arithmetic operations, and solve analogies has significant implications for natural language processing tasks. These properties enable the development of more sophisticated models that can understand and generate human language more effectively. For example, word embeddings can be used as input features for deep learning models in tasks such as sentiment analysis, named entity recognition, and machine translation. They provide a rich representation of words that goes beyond simple one-hot encoding, allowing models to leverage the semantic and syntactic information present in the embeddings. As a result, word embeddings have become a fundamental component in many state-of-the-art natural language processing systems, contributing to significant advancements in the field.

## Visualization Techniques for Word Embeddings and Semantic Spaces

Word embeddings, such as those generated by Word2Vec or GloVe, represent words as high-dimensional vectors in a continuous space. These embeddings capture semantic and syntactic relationships between words, allowing for meaningful analysis and exploration of language. However, visualizing these high-dimensional

vectors can be challenging. One popular technique for visualizing word embeddings is t-Distributed Stochastic Neighbor Embedding (t-SNE).

t-SNE is a non-linear dimensionality reduction algorithm that aims to preserve the local structure of high-dimensional data while mapping it to a lower-dimensional space, typically 2D or 3D. The algorithm works by minimizing the Kullback-Leibler divergence between the probability distributions of pairwise similarities in the high-dimensional and low-dimensional spaces. The resulting visualization reveals clusters of semantically related words, allowing for intuitive exploration of the semantic space.

To apply t-SNE to word embeddings, the first step is to compute the pairwise similarities between words in the high-dimensional space. This can be done using cosine similarity, which measures the cosine of the angle between two vectors. The similarity matrix is then used as input to the t-SNE algorithm, which iteratively adjusts the positions of the words in the low-dimensional space to minimize the divergence between the probability distributions. The optimization process involves a gradient descent algorithm with a carefully chosen learning rate and perplexity parameter, which controls the balance between local and global structure preservation.

Once the t-SNE algorithm converges, the resulting visualization can be plotted using a scatter plot, where each point represents a word, and the proximity of points indicates semantic similarity. The visualization often reveals meaningful clusters of related words, such as synonyms, antonyms, or words belonging to the same semantic category. For example, words related to emotions, such as "happy," "sad," "angry," and "excited," may form distinct clusters in the t-SNE plot. Similarly, words related to animals, such as "dog," "cat," "horse," and "elephant," may group together.

Exploring the semantic space using t-SNE visualizations can provide valuable insights into the structure and organization of language. Researchers and practitioners in natural language processing and computational linguistics can use these visualizations to study semantic relationships, detect patterns, and identify potential biases or inconsistencies in the word embeddings. Additionally, t-SNE visualizations can be used to evaluate the quality of word embeddings by assessing the coherence and interpretability of the resulting clusters. However, it is important to note that t-SNE is a stochastic algorithm, and the exact positions of words in the visualization may vary across different runs. Therefore, multiple visualizations should be generated and compared to ensure the robustness of the observed patterns.

## Beyond Word2Vec

Word2Vec, introduced by Mikolov et al. in 2013, revolutionized the field of natural language processing by providing a simple yet effective method for learning word embeddings. However, since its inception, several modern embedding approaches have emerged, addressing some of the limitations of Word2Vec and extending its capabilities. GloVe (Global Vectors for Word Representation) is one such approach, proposed by Pennington et al. in 2014. GloVe combines the advantages of global matrix factorization and local context window methods, capturing both global word co-occurrence statistics and local context information. The objective function of GloVe is based on the weighted least squares regression model, which aims to minimize the difference between the dot product of word vectors and the logarithm of their co-occurrence probability. The resulting embeddings have been shown to outperform Word2Vec on various word similarity and analogy tasks.

Another notable advancement in word embeddings is the incorporation of subword information, as exemplified by the FastText model introduced by Bojanowski et al. in 2017. FastText extends the Word2Vec model by representing words as a bag of character n-grams, allowing the model to capture morphological information and generate embeddings for out-of-vocabulary words. By considering subword units, FastText can handle rare and unseen words more effectively, making it particularly useful for morphologically rich languages and domains with specialized vocabulary. The model is trained using a modified version of the skip-gram architecture, where the objective is to predict the surrounding words given the character n-grams of the target word. FastText has been shown to improve performance on various downstream tasks, such as text classification and named entity recognition.

A significant paradigm shift in word embeddings occurred with the introduction of contextual embeddings,

which differ from static embeddings like Word2Vec and GloVe. Contextual embeddings, such as those obtained from language models like ELMo (Embeddings from Language Models) and BERT (Bidirectional Encoder Representations from Transformers), generate word representations that are sensitive to the surrounding context. In these models, each word has multiple embeddings depending on its context, allowing for a more nuanced and context-aware representation of words. ELMo, proposed by Peters et al. in 2018, uses a bidirectional LSTM language model to generate contextualized word embeddings. BERT, introduced by Devlin et al. in 2019, employs a transformer architecture with a masked language modeling objective, enabling the model to capture bidirectional context. Contextual embeddings have achieved state-of-the-art performance on a wide range of natural language understanding tasks, such as question answering, sentiment analysis, and named entity recognition.

Moving beyond text-only embeddings, multi-modal embeddings have emerged as a powerful approach to capture cross-modal relationships between different modalities, such as images and text. Multi-modal embeddings aim to learn a joint embedding space where semantically similar concepts from different modalities are mapped close to each other. One prominent example is the Visual-Semantic Embedding (VSE) model, proposed by Frome et al. in 2013, which learns to align visual and textual representations. The VSE model consists of two neural networks: a visual encoder (e.g., a convolutional neural network) that maps images to a visual embedding space, and a text encoder (e.g., a recurrent neural network) that maps sentences to a textual embedding space. The objective is to minimize the distance between the visual and textual embeddings of semantically related pairs while maximizing the distance between unrelated pairs. This allows for cross-modal retrieval, where images can be retrieved based on textual queries and vice versa.

Multi-modal embeddings have found applications in various domains, such as image captioning, visual question answering, and cross-modal retrieval. For example, in image captioning, multi-modal embeddings can be used to generate textual descriptions of images by aligning the visual features with the corresponding word embeddings. In visual question answering, multi-modal embeddings enable the model to reason about the relationship between the question text and the visual content of the image to produce accurate answers. Cross-modal retrieval systems can benefit from multi-modal embeddings by enabling efficient search and retrieval of relevant images based on textual queries or retrieving relevant text documents based on image queries. The success of multi-modal embeddings has led to the development of more advanced models, such as the Transformer-based ViLBERT (Vision-and-Language BERT) and LXMERT (Learning Cross-Modality Encoder Representations from Transformers), which leverage the power of self-attention mechanisms to capture fine-grained interactions between visual and textual features.

## Training Embeddings

The embedding layer is a crucial component in deep learning architectures designed for processing sequential data, such as text or time series. It maps discrete input tokens to dense, continuous vector representations, enabling the model to capture semantic and syntactic relationships between the tokens. The embedding layer is typically implemented as a lookup table, where each row corresponds to a specific token in the vocabulary, and the columns represent the dimensions of the embedding space. During the forward pass, the input tokens are used as indices to retrieve their corresponding embeddings from the lookup table, which are then passed to the subsequent layers of the network.

When training embeddings, the choice of loss function plays a significant role in determining the quality and properties of the learned representations. One commonly used loss function is the cross-entropy loss, which measures the dissimilarity between the predicted probability distribution and the true distribution of the target tokens. This loss function encourages the model to assign higher probabilities to the correct target tokens while minimizing the probabilities of incorrect ones. Another popular choice is the contrastive loss, which aims to bring embeddings of similar tokens closer together in the embedding space while pushing dissimilar tokens apart. This loss function is particularly useful for tasks such as word similarity or analogy completion.

Batch construction is an important consideration when training embeddings, as it directly impacts the efficiency and effectiveness of the learning process. In the context of natural language processing, batches are typically constructed by grouping together sequences of similar lengths to minimize the amount of

padding required. Padding is necessary to ensure that all sequences within a batch have the same length, enabling efficient parallelization on hardware accelerators like GPUs. However, excessive padding can lead to computational waste and potentially degrade the quality of the learned embeddings. To mitigate this issue, techniques such as bucketing can be employed, where sequences are grouped into buckets based on their lengths, and batches are constructed within each bucket.

Another aspect of batch construction is the selection of negative samples for contrastive learning. Negative sampling involves choosing a set of tokens that are dissimilar to the target token, and the model is trained to distinguish between the positive (target) and negative samples. The choice of negative samples can significantly impact the quality of the learned embeddings, as it determines the contrast between similar and dissimilar tokens. Common strategies for negative sampling include random sampling from the vocabulary, frequency-based sampling, and more advanced techniques like noise-contrastive estimation (NCE) or hierarchical softmax.

In addition to the aforementioned considerations, the architecture of the embedding layer itself can influence the training process and the resulting embeddings. One important factor is the dimensionality of the embedding space, which determines the capacity of the embeddings to capture complex relationships between tokens. Higher-dimensional embeddings can potentially capture more nuanced information but may also be more prone to overfitting and require more computational resources. Regularization techniques, such as L2 regularization or dropout, can be applied to the embedding layer to prevent overfitting and improve generalization. Furthermore, the initialization of the embedding weights can impact the convergence speed and the quality of the learned representations. Common initialization strategies include random initialization, pre-trained embeddings (e.g., Word2Vec or GloVe), or task-specific initializations based on domain knowledge.

## Training Strategies for Beyond Word2Vec

Word embeddings have revolutionized the field of natural language processing (NLP) by providing a dense, continuous representation of words in a vector space. While Word2Vec has been a seminal model for learning word embeddings, recent advancements have led to the development of more sophisticated training strategies that go beyond the traditional Word2Vec approach. These strategies aim to capture more complex linguistic properties and improve the performance of downstream NLP tasks.

One such strategy is the distinction between pre-training and task-specific training. Pre-training involves learning general-purpose word embeddings on a large corpus of unlabeled text data, such as Wikipedia or Common Crawl. The objective is to capture the semantic and syntactic relationships between words in an unsupervised manner. Pre-trained embeddings serve as a good starting point for various NLP tasks, as they encapsulate general language knowledge. On the other hand, task-specific training focuses on learning embeddings tailored to a specific task, such as sentiment analysis or named entity recognition. Task-specific embeddings are typically learned by fine-tuning pre-trained embeddings on a labeled dataset relevant to the task at hand. This allows the embeddings to adapt to the specific characteristics and requirements of the target task.

Fine-tuning embeddings is another important training strategy that has gained significant attention in recent years. Instead of using pre-trained embeddings as fixed representations, fine-tuning allows the embeddings to be updated during the training process of a downstream task. By backpropagating the gradients through the embedding layer, the model can adjust the word representations to better suit the task-specific objectives. Fine-tuning has been shown to improve the performance of various NLP tasks, such as text classification, question answering, and machine translation. The process of fine-tuning involves initializing the embedding layer with pre-trained embeddings and then updating them along with the other model parameters during the training phase. This allows the embeddings to capture task-specific nuances and adapt to the domain-specific vocabulary and semantics.

Transfer learning with embeddings is a powerful technique that leverages the knowledge learned from one task to improve the performance on another related task. The idea is to transfer the learned representations from a source task, where abundant labeled data is available, to a target task with limited labeled data. This is particularly useful in scenarios where obtaining large amounts of labeled data for the target task is expensive or time-consuming. Transfer learning with embeddings can be achieved by using pre-trained embeddings

as feature extractors and training a separate model on top of them for the target task. Alternatively, the pre-trained embeddings can be fine-tuned on the target task, allowing them to adapt to the specific characteristics of the new domain. Transfer learning has been successfully applied to various NLP tasks, such as sentiment analysis, named entity recognition, and text classification, demonstrating its effectiveness in leveraging knowledge across different domains and tasks.

The choice of training strategy depends on several factors, including the availability of labeled data, the complexity of the target task, and the computational resources at hand. Pre-training embeddings on large unlabeled corpora can be computationally expensive but provides a good starting point for various downstream tasks. Fine-tuning embeddings allows for task-specific adaptation but requires labeled data and additional training time. Transfer learning offers a way to leverage knowledge from related tasks and can be particularly beneficial when labeled data is scarce for the target task. The optimal training strategy often involves a combination of these approaches, depending on the specific requirements and constraints of the project. Researchers and practitioners need to carefully consider the trade-offs and choose the most suitable strategy based on their specific use case and available resources.

## Handling Challenges in Embeddings

Embeddings, a fundamental component of modern natural language processing (NLP) systems, face several challenges that can impact their effectiveness and generalizability. One such challenge is the handling of rare words, which occur infrequently in the training corpus. Traditional embedding techniques, such as word2vec and GloVe, rely on the co-occurrence statistics of words to learn their vector representations. However, for rare words, the limited number of occurrences may result in poorly learned embeddings that fail to capture the true semantic meaning of these words. This can lead to suboptimal performance in downstream NLP tasks, particularly when the rare words carry significant importance in the context of the application.

Another challenge in embeddings is the presence of out-of-vocabulary (OOV) words. OOV words are those that were not encountered during the training phase and, therefore, do not have a corresponding embedding vector. This issue arises when the embedding model is applied to new or unseen text data that contains words not present in the training vocabulary. Handling OOV words is crucial to ensure the robustness and adaptability of NLP systems in real-world scenarios. Various approaches have been proposed to address this challenge, such as using subword-level embeddings, character-level embeddings, or employing techniques like fastText that leverage morphological information to generate embeddings for OOV words.

Domain adaptation is another significant challenge in embeddings. Embeddings trained on a general corpus may not effectively capture the nuances and domain-specific semantics when applied to specialized domains, such as medical text, legal documents, or technical literature. The vocabulary, writing style, and semantic relationships in these domains often differ from those in general language, leading to a mismatch between the pre-trained embeddings and the target domain. This mismatch can result in suboptimal performance and limited generalization ability of NLP models in domain-specific tasks. To address this challenge, techniques such as domain-specific fine-tuning, transfer learning, and domain adaptation strategies have been proposed. These approaches aim to adapt the pre-trained embeddings to the target domain by leveraging domain-specific corpora or incorporating domain knowledge into the embedding learning process.

The challenges of rare words, out-of-vocabulary words, and domain adaptation in embeddings have been the focus of active research in the NLP community. Various solutions and techniques have been proposed to mitigate these challenges and improve the robustness and generalizability of embedding-based NLP systems. For example, subword-level embeddings, such as byte-pair encoding (BPE) and WordPiece, have been introduced to handle rare and OOV words by breaking them down into smaller, more frequent subword units. These subword units are then used to compose the embeddings of the original words, enabling the model to generate meaningful representations even for unseen or infrequent words. Additionally, character-level embeddings have been explored as an alternative approach, where each word is represented as a sequence of characters, and the embeddings are learned at the character level. This allows the model to capture morphological patterns and handle OOV words by leveraging the compositional nature of language.

To address the challenge of domain adaptation, several techniques have been proposed. One approach is to fine-tune pre-trained embeddings on domain-specific corpora, allowing the embeddings to adapt to the

characteristics and semantics of the target domain. This fine-tuning process involves updating the embedding vectors using the domain-specific data, often in conjunction with the downstream NLP task. Another approach is to employ transfer learning techniques, where embeddings trained on a large, general corpus are used as initialization for domain-specific models. The models are then further trained on the target domain data, leveraging the knowledge captured in the pre-trained embeddings while adapting to the specific domain. Additionally, domain adaptation strategies, such as adversarial training and domain-invariant feature learning, have been explored to learn embeddings that are more robust and transferable across different domains. These strategies aim to align the embedding spaces of the source and target domains, reducing the discrepancy between them and improving the generalization ability of the embeddings.

## Practical Applications of Embeddings in Data Science

Embeddings have found extensive practical applications in various domains of data science, particularly in recommendation systems, information retrieval, and analysis tools. In recommendation systems, embeddings play a crucial role in capturing the latent representations of users and items. User-item embeddings are learned from historical interaction data, such as ratings, clicks, or purchases, to uncover the underlying preferences and similarities. These embeddings enable collaborative filtering techniques, where the preferences of similar users are leveraged to generate personalized recommendations. However, recommendation systems often face the cold start problem, where limited information is available for new users or items. Embeddings can help alleviate this issue by incorporating auxiliary information, such as user demographics or item metadata, to infer initial representations and provide meaningful recommendations.

Embeddings have also revolutionized the field of information retrieval, enabling more effective and efficient search capabilities. By representing documents and queries as dense vectors in a shared embedding space, the similarity between them can be easily computed using vector operations. This allows for semantic search, where the system can retrieve documents that are semantically similar to the query, even if they do not contain exact keyword matches. Embeddings capture the contextual meaning of words and phrases, enabling the retrieval of relevant documents based on their underlying semantics. Furthermore, embeddings have opened up possibilities for cross-lingual information retrieval, where queries in one language can be matched with documents in another language by aligning their embedding spaces. This is achieved through techniques such as cross-lingual word embeddings or multilingual sentence embeddings, which learn a shared representation space across different languages.

In addition to recommendation systems and information retrieval, embeddings serve as powerful tools for various analysis tasks in data science. Bias detection is one such application, where embeddings can be used to uncover and quantify biases present in text data. By analyzing the geometric relationships between word embeddings, researchers can identify and measure gender biases, racial biases, or other forms of discrimination embedded in the language models. This is crucial for developing fair and unbiased AI systems. Embedding probing tasks, on the other hand, aim to understand and interpret the information captured by pre-trained embeddings. By designing specific probing tasks, such as word analogy tests or semantic similarity assessments, researchers can evaluate the quality and properties of embeddings, gaining insights into their strengths and limitations.

The quality evaluation of embeddings is another important aspect in practical applications. Embeddings are often trained on large-scale datasets and can inherit biases or exhibit inconsistencies. Therefore, it is crucial to assess the quality of embeddings before deploying them in real-world systems. Various evaluation metrics and benchmarks have been proposed to measure the semantic coherence, linguistic regularities, and downstream task performance of embeddings. These evaluation techniques help in selecting the most suitable embedding models for specific applications and ensuring their reliability and robustness.

Embeddings have become an indispensable tool in data science, enabling powerful applications in recommendation systems, information retrieval, and analysis tasks. They provide a way to capture and leverage the semantic relationships between entities, whether they are users, items, documents, or words. By representing complex data in a dense and meaningful vector space, embeddings facilitate efficient similarity computations, personalized recommendations, semantic search, and cross-lingual retrieval. Moreover, embeddings serve as a foundation for bias detection, model interpretation, and quality evaluation, ensuring the fairness and reliability

of AI systems. As the field of data science continues to evolve, the practical applications of embeddings are expected to expand further, unlocking new possibilities for intelligent and data-driven solutions.

## Learning Objectives

- Understand the transition from autoencoders to embeddings
- Master word embedding concepts and training
- Implement embedding-based applications
- Evaluate embedding quality and characteristics