

Roisenzvit, Alfredo B.

Working Paper

From Euclidean distance to spatial classification: Unraveling the technology behind GPT models

Serie Documentos de Trabajo, No. 853

Provided in Cooperation with:

University of CEMA, Buenos Aires

Suggested Citation: Roisenzvit, Alfredo B. (2023) : From Euclidean distance to spatial classification: Unraveling the technology behind GPT models, Serie Documentos de Trabajo, No. 853, Universidad del Centro de Estudios Macroeconómicos de Argentina (UCEMA), Buenos Aires

This Version is available at:

<https://hdl.handle.net/10419/297782>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.

UNIVERSIDAD DEL CEMA
Buenos Aires
Argentina

Serie
DOCUMENTOS DE TRABAJO

Área: Ingeniería

**FROM EUCLIDEAN DISTANCE TO SPATIAL CLASSIFICATION:
UNRAVELING THE TECHNOLOGY BEHIND GPT MODELS**

Alfredo B. Roisenzvit

Abril 2023
Nro. 853

https://ucema.edu.ar/publicaciones/doc_trabajo.php
UCEMA: Av. Córdoba 374, C1054AAP Buenos Aires, Argentina
ISSN 1668-4575 (impreso), ISSN 1668-4583 (en línea)
Editor: Jorge M. Streb; asistente editorial: Valeria Dowding <ved@ucema.edu.ar>

From Euclidean Distance to Spatial Classification: Unraveling the Technology behind GPT Models

Alfredo B. Roisenzvit*

Abstract

In this paper, we present a comprehensive analysis of the technology underpinning Generative Pre-trained Transformer (GPT) models, with a particular emphasis on the interrelationships between Euclidean distance, spatial classification, and the functioning of GPT models. Our investigation begins with a thorough examination of Euclidean distance, elucidating its role as a fundamental metric for quantifying the proximity between points in a multi-dimensional space. Following this, we provide an overview of spatial classification techniques, explicating their utility in discerning patterns and relationships within complex data structures. With this foundation, we delve into the inner workings of GPT models, outlining their architectural components, such as the self-attention mechanism and positional encoding. We then explore the process of training GPT models, detailing the significance of tokenization and embeddings. Additionally, we scrutinize the role of Euclidean distance and spatial classification in enabling GPT models to effectively process input sequences and generate coherent output in a wide array of natural language processing tasks. Ultimately, this paper aims to provide a comprehensive understanding of the intricate connections between Euclidean distance, spatial classification, and GPT models, fostering a deeper appreciation of their collective impact on the advancements in artificial intelligence and natural language processing.

* The viewpoints of the author do not necessarily represent the position of the Universidad del CEMA

Introduction

This paper was, in a relevant part, written through prompts fed to GPT-4. The outcomes of such prompts were then analyzed by the human author and ordered into the pre-defined sections of the paper. Although some corrections were made, and the bibliography was properly checked, the objective of the paper was to produce an academic output, via the responses to the mentioned prompts. This, however, was the result of specific questions orderly submitted in prompt form, and a scientific knowledge of the prompted elements was required to generate a conceptually sound paper. As a result, I found a tremendous productivity boost in the process.

Generative Pre-trained Transformers (GPT) are state-of-the-art language models that have demonstrated exceptional capabilities in natural language processing tasks. Understanding the technology behind these models requires a thorough comprehension of the concepts of Euclidean distance and spatial classification, which form the foundation for their architecture.

Euclidean Distance

Euclidean distance, also known as L2 distance or L2 norm, is a fundamental concept in geometry and linear algebra. It is a measure of the straight-line distance between two points in an n-dimensional Euclidean space (Bishop, 2006). Mathematically, the Euclidean distance between two points $A(x_1, x_2, \dots, x_n)$ and $B(y_1, y_2, \dots, y_n)$ can be calculated using the following equation:

$$d(A, B) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$$

In machine learning and data science, Euclidean distance is commonly employed as a similarity measure between data points in tasks such as clustering, classification, and dimensionality reduction (Pedregosa et al., 2011).

Spatial Classification

Spatial classification is a machine learning technique used to categorize data points based on their positions in a multi-dimensional space. It involves learning a decision boundary that separates different classes of data points in the feature space. Distance metrics, such as Euclidean distance, play a crucial role in determining the similarity or dissimilarity between data points and forming the decision boundaries (Bishop, 2006). The logic follows, in a multi-dimensional space, that if the distance between n coordinates on a vector in such space is minimized, then those elements must be similar or equal if such distance is zero or limit thereof. In spatial classification, any given object, real or imaginary, can be given n attributes, and these attributes can then be vectorized in a category pertaining to the object. So, an object can have as much attributes as different classifications we want to use, and then vectorize these attributes with a coordinate system, that will allow therefore the spatial classification techniques.

Common spatial classification techniques include:

3.1 k-Nearest Neighbors (k-NN) k-NN is a non-parametric classification method that considers the k nearest data points to a query point to determine its class (Altman, 1992). The class membership is determined by majority vote, with the query point being assigned to the class most common among its k nearest neighbors.

3.2 Support Vector Machines (SVM) SVM is a linear classification technique that seeks to find the optimal hyperplane separating data points of different classes (Cortes & Vapnik, 1995). The goal is to maximize the margin between the hyperplane and the closest data points from each class, called support vectors. SVM can also be extended to non-linear classification using kernel methods.

3.3 Decision Trees Decision trees are a type of classification model that recursively splits the feature space into regions based on feature values, forming a tree-like structure (Breiman et al., 1984). Each internal node in the tree represents a feature, and each leaf node represents a class label.

GPT Models

GPT models are a type of deep learning architecture that leverages the Transformer model, which was introduced by Vaswani et al. (2017) as a novel approach to sequence-to-sequence learning. The primary components of the Transformer are the self-attention mechanism and the positional encoding, which together enable the model to effectively capture long-range dependencies in the input data.

The Transformer Model: Architecture and Applications

The Transformer model is a groundbreaking architecture designed for sequence-to-sequence tasks in natural language processing and other domains. It has garnered significant attention due to its ability to handle long-range dependencies in input sequences and its highly parallelizable structure. The key components of the Transformer model are the self-attention mechanism and positional encoding, which will be discussed in the context of GPT models.

The Transformer model can be thought of as a highly intelligent language processing machine designed to understand and generate text. Imagine you have a series of jigsaw puzzle pieces (words in a sentence) that you need to put together in the right order to create a

meaningful picture (a coherent sentence). The Transformer model excels at this task, thanks to two key components: the self-attention mechanism and positional encoding.

Let's break down these concepts using real-world metaphoric examples:

1. **Self-Attention Mechanism:** Think of the self-attention mechanism as a group of detectives working together to solve a case. Each detective (word) in the group is responsible for gathering clues about the other detectives and understanding their roles in the case (the sentence). The self-attention mechanism allows the model to figure out which words are closely related and should be considered together to make sense of the whole sentence.

For example, in the sentence, "The dog chased the cat, but it got away," the self-attention mechanism helps the model understand that "dog" is related to "chased" and "cat" is related to "it got away."

2. **Positional Encoding:** The order in which words appear in a sentence often plays a crucial role in determining the meaning. Positional encoding is like a GPS system that gives each word a specific location in the sentence, helping the model understand the importance of the word's position.

For instance, consider the sentences "The cat sat on the mat" and "The mat sat on the cat." Although the words are the same, the meaning of these sentences is different because of the order of the words. Positional encoding helps the model capture the correct meaning by considering the order of the words in the sentence.

Now let's look at the overall structure of the Transformer model. It has two main parts: the encoder and the decoder.

1. Encoder: The encoder's job is to read and analyze the input text, just like a detective gathering information about a case. It takes the words in the sentence, considers their relationships and positions, and creates a continuous representation that captures the meaning of the whole sentence.
2. Decoder: The decoder's role is similar to a detective presenting their findings to solve the case. It takes the continuous representation generated by the encoder, processes it further, and produces an output that can be a translation, a summary, or answers to specific questions.

In summary, the Transformer model is a powerful language processing tool that excels at understanding and generating text by considering the relationships between words and their positions in a sentence. It's like a team of detectives working together to solve complex language puzzles and present meaningful solutions.

Drilling down into the Encoder

In the Transformer model, the encoder plays a crucial role in processing and understanding the input text. Let's break down its workings in layman's terms, using the metaphor of a detective agency and referring to the code as described by Vaswani et al. in the original paper.

The encoder consists of multiple layers, each working like a team of detectives who specialize in different aspects of the case (input text). Each layer in the encoder contains two sub-layers: a multi-head self-attention mechanism and a position-wise feed-forward network.

1. Multi-Head Self-Attention Mechanism: This sub-layer can be thought of as several groups of detectives working simultaneously to analyze different aspects of the case.

Each group (head) focuses on different relationships between the words in the sentence.

By having multiple heads, the model can capture a more comprehensive understanding of the input text.

For example, one group of detectives might focus on understanding the relationships between nouns and verbs, while another group might concentrate on the connections between adjectives and the nouns they describe.

In the original paper, the multi-head self-attention mechanism is defined as:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \text{ where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

Here, Q , K , and V represent query, key, and value matrices. W_i^Q , W_i^K , and W_i^V are the learned linear projections for each head. W^O is the output linear projection.

2. Position-wise Feed-Forward Network: After the self-attention mechanism, the input text is passed through a position-wise feed-forward network, which can be thought of as a group of detectives working together to synthesize the findings of the previous groups. They process the information gathered by the self-attention mechanism and generate a new representation of the input text.

In the original paper, the position-wise feed-forward network is defined as:

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

Here, W_1 and W_2 are weight matrices, and b_1 and b_2 are bias terms. The max function is the Rectified Linear Unit (ReLU) activation function.

The output of the top layer in the encoder is a continuous representation of the input text, which captures the meaning of the whole sentence. This representation is then passed to the decoder to generate the final output.

Throughout the encoder, residual connections and layer normalization are used to improve the training stability and performance. Residual connections can be thought of as a shortcut that allows the detectives to refer back to the original input, while layer normalization helps maintain consistency in the detectives' work, ensuring that their findings are reliable and accurate.

In summary, the encoder in the Transformer model works like a detective agency with multiple layers of specialists who analyze the relationships between words and their positions in a sentence. Through the multi-head self-attention mechanism and position-wise feed-forward networks, the encoder creates a continuous representation of the input text, which is then used by the decoder to generate the final output.

Tokenization

Before the input text is fed into the Transformer model, it needs to be broken down into smaller pieces called tokens. This process is known as tokenization. Let's explain tokenization using a metaphor and in layman's terms, similar to the previous explanations.

Think of tokenization as a librarian who's tasked with organizing books (input text) on shelves in a way that makes it easier for the Transformer model (the reader) to understand and process the information. The librarian breaks down the books into smaller units (tokens), which can be words, subwords, or even characters, depending on the chosen tokenization strategy.

There are several strategies the librarian (tokenizer) can use to break down the text into tokens:

1. **Word-based Tokenization:** In this approach, the librarian breaks down the text into individual words. For example, the sentence "The cat is sitting on the mat" would be

tokenized into ["The", "cat", "is", "sitting", "on", "the", "mat"]. This is a simple and intuitive method, but it may struggle with handling rare or out-of-vocabulary words.

2. Subword-based Tokenization: To overcome the limitations of word-based tokenization, the librarian can use subword-based tokenization, which breaks down the text into smaller units that are more frequent and consistent across languages. For example, the sentence "ChatGPT is an advanced AI model" might be tokenized into ["Chat", "G", "PT", " is", " an", " ad", "van", "ced", " AI", " mo", "del"]. This method can handle rare words and is more flexible when dealing with different languages.
3. Character-based Tokenization: In this approach, the librarian breaks down the text into individual characters. For example, the sentence "The cat is sitting on the mat" would be tokenized into ["T", "h", "e", " ", "c", "a", "t", " ", "i", "s", " ", "s", "i", "t", "t", "i", "n", "g", " ", "o", "n", " ", "t", "h", "e", " ", "m", "a", "t"]. This method can handle any input text but may require more computational resources due to the increased number of tokens.

Once the text is tokenized, the tokens are converted into numerical representations called embeddings. These embeddings are like the Dewey Decimal System used in libraries to represent the tokens in a way that makes it easier for the Transformer model to understand and process the information.

In summary, tokenization is the process of breaking down the input text into smaller units, such as words, subwords, or characters, to make it easier for the Transformer model to understand and process the information. The choice of tokenization strategy depends on the specific problem and desired balance between simplicity, flexibility, and computational resources.

Embeddings and similarities

Embeddings are an essential part of natural language processing, as they help the Transformer model (the reader) understand and process the tokens (pieces of text) more effectively. To explain embeddings using a metaphor and in layman's terms, let's consider the process of assigning colors to objects based on their properties.

Imagine that you have a collection of objects (tokens), such as fruits, and you want to represent their properties, like taste, size, and texture, in a way that's easier for the Transformer model to understand. To do this, you can assign a unique color (embedding) to each fruit, based on its properties. For example, sweet fruits might be represented by shades of red, while sour fruits could be shades of green.

Embeddings work similarly in natural language processing. They are numerical representations that capture the meaning and properties of tokens in a continuous vector space. Instead of representing tokens as discrete units (like colors in our fruit example), embeddings position the tokens in a multi-dimensional space based on their similarities, relationships, and contextual information.

For example, words with similar meanings, such as "cat" and "kitten," would have embeddings that are closer together in the vector space, while unrelated words, like "cat" and "airplane," would be farther apart. This continuous representation allows the Transformer model to process and understand the tokens more effectively.

There are several methods to generate embeddings, including:

1. Word2Vec: This method learns embeddings based on the idea that words that appear in similar contexts tend to have similar meanings. It uses a shallow neural network to

predict a word based on its surrounding words (Skip-gram) or predict surrounding words based on a given word (Continuous Bag of Words).

2. GloVe: The Global Vectors for Word Representation (GloVe) method learns embeddings by analyzing word co-occurrence statistics in a large text corpus. It aims to create embeddings that capture both the local context information (like Word2Vec) and global word co-occurrence patterns.
3. FastText: This method is an extension of Word2Vec that generates embeddings for subword units, such as character n-grams, in addition to whole words. This approach allows FastText to handle rare and out-of-vocabulary words more effectively.
4. BERT: Bidirectional Encoder Representations from Transformers (BERT) is a Transformer-based model that learns contextual embeddings by pre-training on a large text corpus using masked language modeling and next sentence prediction tasks.

Once the embeddings are generated, they serve as the input for the Transformer model, which then processes and understands the text based on the relationships and contextual information captured in these continuous vector representations.

In summary, embeddings are numerical representations of tokens that capture their meaning and properties in a continuous vector space. They help the Transformer model understand and process the input text more effectively by leveraging the similarities, relationships, and contextual information between tokens.

Leveraging similarities in embeddings is crucial for the Transformer model to effectively process and understand the input text. The main idea is that tokens with similar meanings or properties should have embeddings that are closer together in the continuous vector space. By

positioning related tokens close to each other in this space, the model can easily identify their relationships and extract meaningful patterns from the data.

Let's dive deeper into the concept of leveraging similarities in embeddings:

1. **Semantic Similarity:** The similarity between embeddings helps the model capture semantic relationships between tokens. For example, synonyms like "happy" and "joyful" should have embeddings that are close together in the vector space, as they have similar meanings. By representing tokens in this manner, the model can recognize that these words can often be used interchangeably in a sentence and understand their semantic connection.
2. **Syntactic Similarity:** Embeddings can also capture syntactic relationships between tokens. For example, different verb forms such as "run," "running," "ran," and "runs" should have similar embeddings, as they all convey the same action but in different grammatical forms. By capturing syntactic similarities, the model can understand the underlying structure and grammar of the input text, helping it generate coherent and grammatically correct output.
3. **Analogical Reasoning:** Another advantage of leveraging similarities in embeddings is the ability to perform analogical reasoning. For instance, given the relationship "king is to man as queen is to woman," the model can understand the analogy by identifying the relationship between the embeddings of "king" and "man" and applying the same relationship to the embeddings of "queen" and "woman." This ability helps the model make inferences and predictions based on the observed patterns and relationships in the data.

4. Contextual Information: Embeddings can also capture contextual information about how tokens are used in different contexts. For example, the word "bank" can refer to a financial institution or the side of a river, depending on the context. By learning embeddings that reflect these contextual nuances, the model can understand the different meanings of a word based on its surrounding words and generate appropriate output accordingly.

The Transformer model leverages these similarities by processing the input text using self-attention mechanisms and positional encodings. The self-attention mechanism allows the model to focus on different parts of the input sequence, weighting the importance of tokens based on their relationships to other tokens in the sequence. This helps the model capture the contextual information and understand the meaning of tokens in the given context.

In summary, leveraging similarities in embeddings allows the Transformer model to capture semantic and syntactic relationships, perform analogical reasoning, and understand contextual information. These capabilities help the model process the input text more effectively, enabling it to generate accurate and coherent output in various natural language processing tasks.

In the context of leveraging similarities in embeddings, Euclidean distance serves as a metric for measuring the proximity between tokens in the continuous vector space. Embeddings capture semantic, syntactic, and contextual information, positioning tokens with similar properties closer together in the vector space. Euclidean distance quantifies the degree of similarity between these embeddings, allowing the Transformer model to understand and process the relationships between tokens more effectively.

In the context of word embeddings, each token is represented as a point in a high-dimensional vector space, and the Euclidean distance between two embeddings serves as an

indicator of their similarity. Smaller Euclidean distances imply that the tokens are more closely related, while larger distances indicate a weaker relationship. By measuring the Euclidean distances between embeddings, the Transformer model can gauge the semantic and syntactic similarities between tokens, as well as their contextual relationships.

It is important to note that Euclidean distance is not the only metric for measuring similarity between embeddings. Other metrics, such as cosine similarity, are also commonly used in natural language processing tasks. Cosine similarity measures the cosine of the angle between two embeddings, with a value of 1 indicating perfect similarity (parallel vectors) and a value of -1 indicating perfect dissimilarity (antiparallel vectors). Cosine similarity is often preferred in some applications, as it is less sensitive to the magnitude of the embeddings and focuses more on their relative orientations.

In summary, Euclidean distance is a metric that quantifies the degree of similarity between embeddings in the continuous vector space. By measuring the distances between embeddings, the Transformer model can leverage the semantic, syntactic, and contextual information encoded in these representations to effectively process and understand the relationships between tokens in the input text.

Self-Attention Mechanism

The self-attention mechanism allows GPT models to weigh the importance of words in a sequence relative to one another. It computes a set of attention scores for each word in the input sequence using query, key, and value vectors, which are derived from the input word embeddings (Vaswani et al., 2017).

Example 1: Consider a simple sentence, "The cat sat on the mat." In this case, the self-attention mechanism would help the model understand the relationships between the words "cat" and "mat," as well as "sat" and "on," as they have strong contextual relationships.

Example 2: In the sentence, "I arrived at the airport after the flight had departed," the self-attention mechanism allows the model to understand the temporal relationship between "arrived" and "departed," despite the distance between the words in the sequence.

Example 3: For the sentence, "Jane went to the store to buy some groceries, but she forgot her wallet at home," the self-attention mechanism helps the model associate "Jane" with "she" and "her wallet," even though they are separated by several words.

Mathematically, the self-attention mechanism is represented by:

$$\text{Attention}(Q, K, V) = \text{Softmax}((QK^T) / \sqrt{d_k})V$$

where Q, K, and V represent the query, key, and value matrices, and d_k is the dimension of the key vectors.

4.2 Positional Encoding

Positional encoding is used to incorporate the order of words in a sequence, as the self-attention mechanism is permutation invariant (Vaswani et al., 2017). It uses sine and cosine functions to encode the position of each word, allowing the model to capture the relative positions of words.

Example 1: Consider the sentences, "I love chocolate ice cream" and "I love ice cream chocolate." Although the words are the same, the meaning of the sentences is different due to the word order. Positional encoding allows the model to differentiate between these sentences by capturing the order of the words.

Example 2: In the sentence, "The quick brown fox jumps over the lazy dog," the word "jumps" is positioned between "fox" and "dog." Positional encoding ensures that the model understands the importance of this position and can capture the relationship between these three words.

Example 3: For the sentence, "Anna bought a new dress, and Emily purchased a stylish skirt," positional encoding enables the model to associate "Anna" with "dress" and "Emily" with "skirt" by capturing the positions of these words in the sequence.

Mathematically, positional encoding is represented by:

$$PE(pos, 2i) = \sin(pos / 10000^{(2i / d_model)})$$

$$PE(pos, 2i + 1) = \cos(pos / 10000^{(2i / d_model)})$$

where pos is the position of the word in the sequence, i is the dimension of the word embedding, and d_model is the model's input dimension.

By integrating self-attention and positional encoding, GPT models can effectively capture the complex relationships between words in a sequence and generate contextually relevant output for various natural language processing tasks.

Training the Model

Training the Transformer model is like teaching a student to understand and generate text based on patterns and relationships learned from large amounts of data. In the context of our metaphors, it involves refining the process of assigning colors to objects (embeddings) and teaching the detective agency (encoder) and the writer (decoder) to work together more efficiently.

Training supervision techniques

For example, GPT-3, or the third iteration of the Generative Pre-trained Transformer, is primarily trained using unsupervised learning techniques with a small amount of supervised fine-tuning. The training process can be broken down into two main stages: pre-training and fine-tuning.

1. Pre-training (unsupervised learning): The primary learning technique employed in GPT-3's training is unsupervised learning. During this stage, the model learns from vast amounts of text data without any explicit supervision or labeled data. The primary objective is to learn a language model that can predict the next word in a sequence based on the given context. This is achieved using a technique called maximum likelihood estimation.

The pre-training process involves feeding the model with a large corpus of text from various sources, such as websites, books, articles, and more. The model tokenizes the text into smaller units (words or subwords) and processes them using the Transformer architecture.

The Transformer architecture consists of multiple layers of self-attention mechanisms and feed-forward neural networks, which help the model learn contextual representations of the input tokens. As the model processes the input, it predicts the next token in the sequence. During pre-training, the model's parameters (weights and biases) are updated using backpropagation to minimize the prediction error. This process continues iteratively until the model converges, meaning its predictions become more accurate and the error is minimized.

2. Fine-tuning (supervised learning): After pre-training, GPT-3 is fine-tuned using a smaller, labeled dataset. This dataset contains input-output pairs, where the input is a prompt and the output is the desired response. Fine-tuning is a form of supervised

learning because the model is provided with explicit supervision in the form of labeled data.

During fine-tuning, the model is trained to minimize the difference between its predictions and the provided target outputs. The model's parameters are further updated using techniques like gradient descent, but the updates are smaller compared to the pre-training phase. This fine-tuning process allows GPT-3 to adapt to specific tasks, such as translation, summarization, or question-answering, depending on the labeled dataset used.

In summary, GPT-3's training process primarily involves unsupervised learning during the pre-training phase, where the model learns to predict the next word in a sequence based on context. After pre-training, the model is fine-tuned on a smaller, task-specific labeled dataset using supervised learning techniques, allowing it to adapt to specific tasks and generate more accurate outputs.

Training steps

Let's break down the training process into three main steps: pre-processing, forward propagation, and backpropagation.

1. Pre-processing: The first step in training the model is preparing the data. The input text is tokenized into smaller units (words, subwords, or characters), and these tokens are converted into embeddings (assigning colors to objects). The target output (e.g., translations, summaries, or labeled entities) is also tokenized and embedded.
2. Forward Propagation: During this step, the input embeddings are passed through the encoder (detective agency) to generate a continuous representation of the input text. This representation captures the relationships and contextual information between the tokens.

The decoder (writer) then uses this representation to generate the output text one token at a time, predicting the most likely next token based on the input, previous output tokens, and its internal state. The model's predictions are compared to the actual target output to calculate a loss, which measures how well the model is performing.

3. Backpropagation: To improve the model's performance, the loss is used to update the model's parameters (the weights and biases in the self-attention mechanism, position-wise feed-forward networks, and other components). The backpropagation algorithm calculates the gradients (partial derivatives) of the loss with respect to each parameter, indicating how much each parameter should be adjusted to minimize the loss. The model's parameters are then updated using an optimization algorithm, such as stochastic gradient descent (SGD) or Adam, which adjusts the parameters based on the calculated gradients and a learning rate.

The training process is repeated for many iterations (epochs) over the entire dataset, with the model's performance improving incrementally as it learns the patterns and relationships in the data. To prevent overfitting, techniques such as dropout, weight decay, and early stopping can be employed, ensuring that the model generalizes well to new, unseen data.

In summary, the training of the Transformer model involves pre-processing the input and target text, passing the data through the encoder and decoder to generate predictions, and using the calculated loss to update the model's parameters. This process is iterated over the dataset, refining the model's understanding of the input text and its ability to generate accurate and coherent output.

Drilling down on how the model predicts the next token: *this is the magic sauce.*

The prediction of the next token involves a series of steps within the Transformer architecture, which is the underlying model for GPT. Here, we'll dive into a technical explanation of how GPT-3 predicts the next token in a given sequence:

1. Tokenization: The input text is first tokenized into smaller units, such as words or subwords. These tokens are then converted into numerical representations called embeddings using an embedding layer. Each token is also assigned a positional encoding to retain information about the position of the token in the sequence.
2. Forward pass through the Transformer layers: GPT-3 is a deep learning model that consists of multiple Transformer layers stacked on top of each other. Each layer is composed of a multi-head self-attention mechanism and a position-wise feed-forward neural network, with layer normalization and residual connections in between.
 - a. Multi-head self-attention: The self-attention mechanism computes a weighted sum of the input embeddings based on their similarity. This process is done in parallel for multiple "heads," which allows the model to capture different aspects of the input relationships. The outputs of all the heads are concatenated and passed through a linear layer.
 - b. Position-wise feed-forward neural network: After the multi-head self-attention step, the output is passed through a feed-forward neural network consisting of two linear layers separated by an activation function, usually ReLU (Rectified Linear Unit).
 - c. Layer normalization and residual connections: Layer normalization is applied after the multi-head self-attention and the feed-forward neural network to stabilize the learning process. Residual connections are used to add the input of each sub-layer to its output, which helps in mitigating the vanishing gradient problem.

3. Repeating the process across layers: The output of one Transformer layer becomes the input for the next layer. This process is repeated for all the layers in the model. The output of the final layer is a contextualized representation of the input tokens.
4. Generating the probability distribution for the next token: The output of the final Transformer layer is passed through a linear layer to project the contextualized embeddings back into the vocabulary size. Then, a softmax activation function is applied, which converts the logits into a probability distribution over the vocabulary. The softmax function ensures that the sum of probabilities for all possible tokens is equal to 1.
5. Sampling or selecting the next token: To predict the next token, one can either choose the token with the highest probability (greedy decoding) or sample a token from the probability distribution (sampling-based decoding). The chosen token is then added to the input sequence, and the process is repeated until a predetermined stopping condition is met (e.g., reaching a maximum sequence length or generating an end-of-sequence token).

In summary, GPT-3 predicts the next token in a sequence by employing a sophisticated multi-layered Transformer architecture, which comprises several key components that work in tandem to process the input sequence effectively. These components include self-attention mechanisms, feed-forward neural networks, layer normalization, and residual connections. Upon processing the input sequence through the multiple layers of the Transformer architecture, GPT-3 generates a final output in the form of a probability distribution over the entire vocabulary. This distribution represents the model's confidence in each potential token being the most suitable next token in the sequence, with higher probabilities assigned to tokens deemed more likely based on the contextual information derived from the input. The chosen token is then appended to the input sequence, and the process is repeated until a predefined stopping criterion

is met, such as reaching a maximum sequence length or generating a special end-of-sequence token. This comprehensive approach, combining self-attention mechanisms, feed-forward neural networks, layer normalization, and residual connections, equips GPT-3 with the ability to generate highly coherent and contextually appropriate output, making it a powerful model for a wide range of natural language processing tasks.+

Conclusion

Although GPT models are primarily designed for natural language processing, the concepts of Euclidean distance and spatial classification can be intimately linked to their underlying mechanisms, highlighting the versatility and adaptability of these models for various practical applications.

The self-attention mechanism, a crucial component of GPT models, computes similarity scores between words in a sequence. This process is analogous to the role of distance metrics, such as Euclidean distance, in spatial classification, where the proximity between data points is used to discern patterns and relationships within complex data structures. Furthermore, GPT models can be adapted for tasks involving spatial data by leveraging spatial classification techniques. These adaptations may include the incorporation of convolutional layers, attention mechanisms designed for spatial data, or specialized positional encoding schemes that capture spatial relationships between elements. Some examples of tasks that can benefit from such adaptations include:

- Image classification: GPT models can be modified to handle image data by incorporating convolutional layers, which can extract local features and spatial patterns in images.

These features can then be processed by the Transformer architecture to generate a classification output.

- **Object detection:** By integrating specialized attention mechanisms that focus on spatial relationships and local contexts, GPT models can be adapted to detect and localize objects within images, providing both class labels and bounding box coordinates.
- **Semantic segmentation:** GPT models can be extended to perform pixel-wise classification, assigning class labels to each pixel in an image. This task can benefit from adaptations that maintain and utilize spatial information throughout the model, such as preserving spatial dimensions in intermediate layers or employing specialized positional encoding schemes.

In conclusion, the connection between Euclidean distance, spatial classification, and GPT models reveals the potential for these models to be adapted and fine-tuned for a wide range of practical applications, transcending their original purpose in natural language processing. By understanding and exploiting these connections, researchers and practitioners can unlock new possibilities and drive innovation in artificial intelligence across various domains.

References

- Altman, N. S. (1992). An introduction to kernel and nearest-neighbor nonparametric regression. *The American Statistician*, 46(3), 175-185.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
- Breiman, L., Friedman, J., Olshen, R. A., & Stone, C. J. (1984). *Classification and Regression Trees*. CRC Press.
- Bojanowski, Piotr, et al. "Enriching Word Vectors with Subword Information." *Transactions of the Association for Computational Linguistics*, vol. 5, 2017, pp. 135-146.
- Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3), 273-297.
- Dai, A. M., Yang, Y., Nguyen, T., Le, Q. V., & Li, F.-F. (2019). Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context. *Association for Computational Linguistics*. <https://doi.org/10.18653/v1/N19-4006>
- Devlin, Jacob, et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 2019, pp. 4171-4186.
- Lubbad, M. The Ultimate Guide to GPT-4 Parameters Web: <https://medium.com/@mlubbad/the-ultimate-guide-to-gpt-4-parameters-everything-you-need-to-know-about-nlps-game-changer-109b8767855a>
- Mikolov, Tomas, et al. "Efficient Estimation of Word Representations in Vector Space." *Proceedings of the International Conference on Learning Representations (ICLR)*, 2013.
- OpenAI. (n.d.). GPT-4. <https://openai.com/gpt-4/>

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.

Pennington, Jeffrey, Richard Socher, and Christopher D. Manning. "GloVe: Global Vectors for Word Representation." *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1532-1543.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998-6008.