

# Feature Group Embeddings: Comprehensive Guide

Aryan Dangi

December 18, 2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Understanding Tags and Tag Dimensions</b>	<b>3</b>
2.1	Tags . . . . .	3
2.2	Tag Dimensions . . . . .	3
<b>3</b>	<b>Why Use Embeddings for Tags?</b>	<b>3</b>
3.1	Challenge with Categorical Data . . . . .	3
3.2	Solution with Embeddings . . . . .	3
3.3	Advantages . . . . .	4
<b>4</b>	<b>Feature Group Embeddings Methods</b>	<b>4</b>
4.1	Sum of Embeddings . . . . .	4
4.1.1	Conceptual Overview . . . . .	4
4.1.2	Detailed Working . . . . .	4
4.1.3	Code Explanation . . . . .	4
4.1.4	Advantages and Considerations . . . . .	7
4.2	Linear Mapping (Concatenation via Learnable Weight Matrix) . . . . .	7
4.2.1	Conceptual Overview . . . . .	7
4.2.2	Detailed Working . . . . .	7
4.2.3	Code Explanation . . . . .	8
4.2.4	Advantages and Considerations . . . . .	8
4.3	Multi-Head Embeddings . . . . .	9
4.3.1	Conceptual Overview . . . . .	9
4.3.2	Detailed Working . . . . .	9
4.3.3	Code Explanation . . . . .	9
4.3.4	Advantages and Considerations . . . . .	10
4.4	Interaction Embeddings . . . . .	10
4.4.1	Conceptual Overview . . . . .	10
4.4.2	Detailed Working . . . . .	11
4.4.3	Code Explanation . . . . .	11
4.4.4	Advantages and Considerations . . . . .	13
4.5	Neural Network-Based Embeddings (MLP Approach) . . . . .	13
4.5.1	Conceptual Overview . . . . .	13
4.5.2	Detailed Working . . . . .	13
4.5.3	Code Explanation . . . . .	14

4.5.4	Advantages and Considerations	14
4.6	Bilinear Embeddings	15
4.6.1	Conceptual Overview	15
4.6.2	Detailed Working	15
4.6.3	Code Explanation	15
4.6.4	Advantages and Considerations	16
4.7	Attention-Based Embeddings	16
4.7.1	Conceptual Overview	16
4.7.2	Detailed Working	16
4.7.3	Code Explanation	17
4.7.4	Correct Implementation for Per-Tag Attention	18
4.7.5	Advantages and Considerations	19
4.8	Hierarchical Embeddings	19
4.8.1	Conceptual Overview	19
4.8.2	Detailed Working	19
4.8.3	Code Explanation	20
4.8.4	Advantages and Considerations	21
4.9	Pretrained Embeddings	21
4.9.1	9.1 Conceptual Overview	21
4.9.2	Detailed Working	22
4.9.3	Code Explanation	22
4.9.4	Advantages and Considerations	23
<b>5</b>	<b>Choosing the Right Embedding Method</b>	<b>23</b>
<b>6</b>	<b>Integrating Embeddings into the TimeXer Model</b>	<b>24</b>
6.1	Example Integration Using Sum of Embeddings	24
6.1.1	Explanation	25
6.1.2	Usage Example	25
<b>7</b>	<b>Final Recommendations and Best Practices</b>	<b>26</b>
<b>8</b>	<b>Conclusion</b>	<b>27</b>

# 1 Introduction

In machine learning, particularly when dealing with complex datasets, it's common to encounter features categorized or grouped based on certain attributes or "tags." Feature Group Embeddings transform these categorical tags into continuous numerical representations, enabling models to process them effectively. This guide explores various embedding methods, their implementations, advantages, and considerations, providing a comprehensive understanding to aid in selecting the most suitable approach for your models.

## 2 Understanding Tags and Tag Dimensions

### 2.1 Tags

- **Definition:** Tags are categorical labels assigned to features to denote their characteristics, categories, or groups. For example, in a dataset containing various sensor readings, tags might represent the type of sensor (temperature, humidity, pressure), the location of the sensor (indoor, outdoor), or the unit of measurement (Celsius, Fahrenheit).
- **Purpose:** Tags help in organizing features, enabling the model to understand relationships and similarities between different features based on their categories.

### 2.2 Tag Dimensions

- **Definition:** The tag dimension refers to the number of unique tags (categories) assigned to features. If there are  $K$  unique tags, the tag dimension is  $K$ .
- **Example:**
  - **SensorType:** Temperature, Humidity, Pressure (3 categories)
  - **Location:** Indoor, Outdoor (2 categories)
  - **Unit:** Celsius, Fahrenheit, Kelvin (3 categories)
  - **Total Unique Tags ( $K$ ):**  $3 + 2 + 3 = 8$  tags

## 3 Why Use Embeddings for Tags?

### 3.1 Challenge with Categorical Data

Machine learning models, especially neural networks, perform best with numerical input. However, categorical data (like tags) are non-numeric and can't be directly fed into these models. Traditional methods like one-hot encoding can lead to high-dimensional and sparse representations, which are inefficient and may not capture the inherent relationships between categories.

### 3.2 Solution with Embeddings

Embeddings transform categorical tags into dense, continuous vectors of lower dimensionality. These embeddings capture semantic relationships between tags, allowing the model to learn similarities and differences effectively.

### 3.3 Advantages

- **Dimensionality Reduction:** Converts high-dimensional categorical data into manageable lower dimensions (e.g., 8 tags  $\rightarrow$  16-dimensional embedding space).
- **Capturing Relationships:** Similar tags can have similar embeddings, enabling the model to generalize better.

## 4 Feature Group Embeddings Methods

There are several methods to embed feature tags, each with its unique advantages and considerations. Below, we explore each method in detail.

### 4.1 Sum of Embeddings

#### 4.1.1 Conceptual Overview

**Sum of Embeddings** involves assigning a unique embedding vector to each tag and summing the embeddings of all active tags for a feature. This method treats each tag independently and aggregates their contributions linearly.

#### 4.1.2 Detailed Working

- **Tag Embedding Assignment:** Each unique tag is associated with a learnable embedding vector of fixed dimension (`embedding_dim`). For  $K$  unique tags, this results in an embedding matrix of size  $(K, \text{embedding\_dim})$ .
- **Feature Representation:** For a feature with multiple tags, the binary tag vector indicates which tags are active. The embeddings corresponding to these active tags are retrieved and summed element-wise to produce the feature's embedding.
- **Mathematical Formulation:**

$$\text{TagEmbed}(f_i) = \sum_{j=1}^K T_i[j] \cdot E_j$$

Where:

- $T_i[j]$  is 1 if tag  $j$  is active for feature  $f_i$ , else 0.
- $E_j$  is the embedding vector for tag  $j$ .

#### 4.1.3 Code Explanation

```
1 import torch
2 import torch.nn as nn
3 import pandas as pd
4 import numpy as np
5
6 # Parameters
7 num_tags = 17
8 embedding_dim = 8
9 output_dim = 1
```

```

10 input_csv = '/kaggle/input/jane-street-real-time-market-data-forecasting
    /features.csv' # Path to your input dataset
11 output_csv = '/kaggle/working/tag_embeddings.csv' # Path to save the
    processed dataset
12
13 # Define the SumEmbedding class
14 class SumEmbedding(nn.Module):
15     def __init__(self, num_tags, embedding_dim):
16         """
17         Initializes the SumEmbedding layer.
18
19         Args:
20             num_tags (int): Total number of unique tags.
21             embedding_dim (int): Dimension of each embedding vector.
22         """
23         super(SumEmbedding, self).__init__()
24         self.tag_embeddings = nn.Embedding(num_tags, embedding_dim)
25         # Initialize embeddings with Xavier Uniform for better
convergence
26         nn.init.xavier_uniform_(self.tag_embeddings.weight)
27
28     def forward(self, tag_vectors):
29         """
30         Computes the sum of embeddings for active tags.
31
32         Args:
33             tag_vectors (torch.Tensor): Binary tensor of shape (
num_features, num_tags)
34
35         Returns:
36             torch.Tensor: Summed embeddings of shape (num_features,
embedding_dim)
37         """
38         embeddings = self.tag_embeddings.weight # (num_tags,
embedding_dim)
39         embedded = torch.matmul(tag_vectors, embeddings) # (
num_features, embedding_dim)
40         return embedded
41
42 def preprocess_data(input_csv, output_csv, num_tags, embedding_dim):
43     """
44     Preprocesses the dataset by applying sum of embeddings to tag
vectors.
45
46     Args:
47         input_csv (str): Path to the input CSV file.
48         output_csv (str): Path to save the processed CSV file.
49         num_tags (int): Number of unique tags.
50         embedding_dim (int): Dimension of the embeddings.
51     """
52     # Load the dataset
53     df = pd.read_csv(input_csv)
54
55     # Identify tag columns
56     tag_columns = [f'tag_{i}' for i in range(num_tags)] # ['tag_0', '
tag_1', ..., 'tag_16']
57
58     # Verify tag columns exist

```

```

59     missing_tags = [tag for tag in tag_columns if tag not in df.columns]
60     if missing_tags:
61         raise ValueError(f"Missing tag columns in the dataset: {
missing_tags}")
62
63     # Initialize the embedding layer
64     sum_embedding = SumEmbedding(num_tags, embedding_dim)
65
66     # Set the embedding layer to evaluation mode to prevent gradient
computation
67     sum_embedding.eval()
68
69     # Extract tag vectors for all features
70     # Shape: (num_features, num_tags)
71     tag_vectors = df[tag_columns].values # numpy array
72
73     # Convert to torch tensor
74     tag_vectors_tensor = torch.tensor(tag_vectors, dtype=torch.float32)
75     # (num_features, num_tags)
76
77     # Compute summed embeddings
78     with torch.no_grad():
79         embeddings = sum_embedding(tag_vectors_tensor) # (num_features,
embedding_dim)
80
81     # Convert embeddings to numpy
82     embeddings_np = embeddings.numpy() # (num_features, embedding_dim)
83
84     # Add embedding columns to the dataframe
85     for dim in range(embedding_dim):
86         embedding_col_name = f'embed_{dim}'
87         df[embedding_col_name] = embeddings_np[:, dim]
88
89     # Optionally, drop original tag columns to reduce data size
90     # Uncomment the following line if you wish to remove tag columns
91     # df.drop(columns=tag_columns, inplace=True)
92
93     # Save the processed dataset
94     df.to_csv(output_csv, index=False)
95     print(f"Processed data saved to {output_csv}")
96
97 preprocess_data(input_csv, output_csv, num_tags, embedding_dim)

```

Listing 1: Sum of Embeddings Implementation

- **Initialization (`__init__`):**

- `num_tags`: Total number of unique tags ( $K$ ).
- `embedding_dim`: Dimension of each embedding vector ( $d_e$ ).
- `self.tag_embeddings`: An embedding layer that maps each tag index to its embedding vector.

- **Forward Pass (`forward`):**

- `tag_vectors`: A binary matrix of shape  $(\text{batch\_size}, \text{num\_tags})$ , where each row represents the active tags for a feature.

- `self.tag_embeddings.weight`: Retrieves the embedding matrix of shape `(num_tags, embedding_dim)`.
- `embedded = tag_vectors @ embeddings`: Performs a matrix multiplication between the binary tag vectors and the embedding matrix, resulting in a summed embedding for each feature in the batch. The resulting shape is `(batch_size, embedding_dim)`.

#### 4.1.4 Advantages and Considerations

##### Advantages:

- **Simplicity**: Easy to implement and computationally efficient.
- **Parameter Efficiency**: Requires only  $K \times d_e$  parameters.
- **Flexibility**: Can handle features with multiple active tags seamlessly.

##### Considerations:

- **Additive Nature**: Assumes each tag contributes independently and additively to the final embedding, potentially missing complex interactions between tags.
- **Independence Assumption**: Does not model dependencies or interactions between tags beyond simple addition.

## 4.2 Linear Mapping (Concatenation via Learnable Weight Matrix)

### 4.2.1 Conceptual Overview

**Linear Mapping** uses a learnable weight matrix to linearly map the binary tag vectors into a continuous embedding space, allowing the model to learn optimal combinations of tags.

### 4.2.2 Detailed Working

- **Weight Matrix ( $W_{\text{tags}}$ )**: A learnable matrix of shape `(embedding_dim, num_tags)`. Each column corresponds to a tag's contribution to each dimension of the embedding.
- **Feature Embedding**: For a given feature's binary tag vector, the embedding is computed as the linear combination of the active tags' columns in  $W_{\text{tags}}$ .
- **Mathematical Formulation**:

$$\text{TagEmbed}(f_i) = W_{\text{tags}} \cdot T_i^T$$

Where:

- $T_i$  is the binary tag vector for feature  $f_i$ .
- $W_{\text{tags}}$  is the weight matrix of shape `(embedding_dim, num_tags)`.

### 4.2.3 Code Explanation

```
1 class LinearEmbedding(nn.Module):
2     def __init__(self, num_tags, embedding_dim):
3         super(LinearEmbedding, self).__init__()
4         self.W_tags = nn.Parameter(torch.randn(embedding_dim, num_tags))
5
6     def forward(self, tag_vectors):
7         """
8         Args:
9             tag_vectors (torch.Tensor): Binary tensor of shape (
batch_size, num_tags)
10
11         Returns:
12             torch.Tensor: Linearly mapped embeddings of shape (
batch_size, embedding_dim)
13         """
14         embedded = tag_vectors.float() @ self.W_tags.t() # (batch_size,
embedding_dim)
15         return embedded
```

Listing 2: Linear Mapping Implementation

- **Initialization (`__init__`):**

- `self.W_tags`: A learnable weight matrix initialized with random values. Shape:  $(\text{embedding\_dim}, \text{num\_tags})$ .

- **Forward Pass (`forward`):**

- `tag_vectors`: Binary matrix of shape  $(\text{batch\_size}, \text{num\_tags})$ .
- `tag_vectors.float()`: Converts the binary tensor to floating-point for matrix multiplication.
- `self.W_tags.t()`: Transposes `W_tags` to shape  $(\text{num\_tags}, \text{embedding\_dim})$  to align for multiplication.
- `embedded = tag_vectors.float() @ self.W_tags.t()`: Multiplies the binary tag vectors with the transposed weight matrix, resulting in  $(\text{batch\_size}, \text{embedding\_dim})$ .

### 4.2.4 Advantages and Considerations

#### Advantages:

- **Learned Combinations**: Enables the model to learn how different tags interact and contribute to the embedding.
- **Expressiveness**: More capable of capturing complex relationships between tags compared to simple summation.

#### Considerations:

- **Increased Parameters**: Requires  $\text{embedding\_dim} \times K$  parameters, which can be substantial for large  $K$ .



- **Risk of Overfitting:** More parameters increase the likelihood of overfitting, especially with limited training data.
- **Computational Cost:** Higher memory and computation requirements compared to simpler methods like summation.

## 4.3 Multi-Head Embeddings

### 4.3.1 Conceptual Overview

**Multi-Head Embeddings** take inspiration from the multi-head attention mechanism in Transformers. Instead of having a single embedding vector per tag, each tag has multiple embedding "heads." This allows capturing different aspects or dimensions of each tag's representation.

### 4.3.2 Detailed Working

- **Multiple Embedding Heads:** Each tag is associated with `num_heads` separate embedding vectors. The final embedding is a concatenation or combination of these heads.
- **Representation Splitting:** The embedding dimension is divided equally among the heads. For example, with `embedding_dim = 64` and `num_heads = 4`, each head has a dimension of 16.
- **Mathematical Formulation:**

$$E_{j,h} \in \mathbb{R}^{\text{head\_dim}}$$

Where  $E_{j,h}$  is the embedding vector for tag  $j$  in head  $h$ .

- The final embedding for a feature is a combination (e.g., concatenation) of embeddings from all heads.

### 4.3.3 Code Explanation

```

1 class MultiHeadEmbedding(nn.Module):
2     def __init__(self, num_tags, embedding_dim, num_heads=4):
3         super(MultiHeadEmbedding, self).__init__()
4         assert embedding_dim % num_heads == 0, "embedding_dim must be
5             divisible by num_heads"
6         self.num_heads = num_heads
7         self.head_dim = embedding_dim // num_heads
8         self.tag_embeddings = nn.Embedding(num_tags, embedding_dim)
9
10        def forward(self, tag_vectors):
11            embeddings = self.tag_embeddings.weight # (num_tags,
12            embedding_dim)
13            embedded = tag_vectors @ embeddings # (batch_size,
14            embedding_dim)
15            # Split into multiple heads
16            embedded = embedded.view(-1, self.num_heads, self.head_dim) # (
17            batch_size, num_heads, head_dim)

```

```
return embedded
```

Listing 3: Multi-Head Embeddings Implementation

- **Initialization (`--init--`):**
  - `num_heads`: Number of embedding heads (default is 4).
  - `self.head_dim`: Dimension of each head, ensuring `embedding_dim` is divisible by `num_heads`.
  - `self.tag_embeddings`: Embedding layer with shape `(num_tags, embedding_dim)`.
- **Forward Pass (`forward`):**
  - `embeddings = self.tag_embeddings.weight`: Retrieves embedding matrix `(num_tags, embedding_dim)`.
  - `embedded = tag_vectors @ embeddings`: Multiplies binary tag vectors with embedding matrix to get `(batch_size, embedding_dim)`.
  - `embedded.view(-1, self.num_heads, self.head_dim)`: Reshapes the embeddings to separate the heads, resulting in `(batch_size, num_heads, head_dim)`.

**Note:** This implementation returns multiple heads separately. Further processing (like attention mechanisms) would typically operate on these heads.

#### 4.3.4 Advantages and Considerations

##### Advantages:

- **Enhanced Representation:** Captures multiple facets or perspectives of each tag’s meaning.
- **Parallel Learning:** Similar to multi-head attention, allowing the model to learn diverse aspects simultaneously.

##### Considerations:

- **Increased Complexity:** More complex architecture requires careful integration with downstream components.
- **Higher Computational Overhead:** More parameters and operations due to multiple heads.
- **Potential Underutilization:** If not properly leveraged, multiple heads might not provide significant benefits over single-head embeddings.

## 4.4 Interaction Embeddings

### 4.4.1 Conceptual Overview

**Interaction Embeddings** aim to explicitly model interactions between pairs (or higher-order combinations) of tags. This is similar to techniques used in factorization machines, where interactions between features are crucial for capturing dependencies.

#### 4.4.2 Detailed Working

- **Pairwise Interactions:** For each pair of active tags in a feature, an interaction term is computed and incorporated into the final embedding.
- **Feature Combination:** The interactions can be modeled by additional parameters or by augmenting the embedding vectors.
- **Mathematical Formulation:**

$$\text{Interaction}_{jk} = E_j \odot E_k$$

Where  $\odot$  denotes element-wise multiplication.

The final embedding could then be:

$$\text{TagEmbed}(f_i) = \sum_{j=1}^K T_i[j] \cdot E_j + \sum_{j=1}^K \sum_{k=j+1}^K T_i[j] \cdot T_i[k] \cdot \text{Interaction}_{jk}$$

#### 4.4.3 Code Explanation

```
1 class InteractionEmbedding(nn.Module):
2     def __init__(self, num_tags, embedding_dim, interaction_order=2):
3         super(InteractionEmbedding, self).__init__()
4         self.num_tags = num_tags
5         self.embedding_dim = embedding_dim
6         self.interaction_order = interaction_order
7         self.tag_embeddings = nn.Embedding(num_tags, embedding_dim)
8
9     def forward(self, tag_vectors):
10        embeddings = self.tag_embeddings.weight # (num_tags,
embedding_dim)
11        embedded = tag_vectors @ embeddings # (batch_size,
embedding_dim)
12
13        if self.interaction_order >= 2:
14            # Pairwise interactions
15            interactions = torch.matmul(tag_vectors, tag_vectors.t()) #
Simplistic example
16            # More sophisticated interaction modeling can be implemented
17            # For illustration, we concatenate
18            embedded = torch.cat([embedded, interactions], dim=1)
19
20        return embedded
```

Listing 4: Interaction Embeddings Implementation

- **Initialization (`__init__`):**
  - `interaction_order`: Determines the order of interactions to model (default is 2 for pairwise).
  - `self.tag_embeddings`: Embedding layer (`num_tags`, `embedding_dim`).
- **Forward Pass (`forward`):**

- `embeddings = self.tag_embeddings.weight:`  
Retrieves `(num_tags, embedding_dim)`.
- `embedded = tag_vectors @ embeddings:`  
Computes summed embeddings `(batch_size, embedding_dim)`.
- **Pairwise Interactions:**
  - \* `interactions = torch.matmul(tag_vectors, tag_vectors.t()):` **Note:**  
This line attempts to compute interactions but is incorrectly implemented for batch-wise pairwise interactions. It incorrectly computes a batch-wise matrix multiplication across batches, resulting in `(batch_size, batch_size)`.
  - \* `embedded = torch.cat([embedded, interactions], dim=1):` Concatenates the original embedding with the interaction terms, increasing the embedding dimension.

**Note:** The provided code has an issue in computing pairwise interactions correctly. Here's a corrected version for pairwise interactions within each sample:

```

1 class InteractionEmbedding(nn.Module):
2     def __init__(self, num_tags, embedding_dim, interaction_order=2):
3         super(InteractionEmbedding, self).__init__()
4         self.num_tags = num_tags
5         self.embedding_dim = embedding_dim
6         self.interaction_order = interaction_order
7         self.tag_embeddings = nn.Embedding(num_tags, embedding_dim)
8
9     def forward(self, tag_vectors):
10        batch_size = tag_vectors.size(0)
11        embedded = tag_vectors @ self.tag_embeddings.weight # (
12        batch_size, embedding_dim)
13
14        if self.interaction_order >= 2:
15            # Compute pairwise interactions for each sample
16            interaction_embeddings = []
17            for i in range(batch_size):
18                active_tags = torch.nonzero(tag_vectors[i], as_tuple=
19                False)[: ,1]
20                if len(active_tags) < 2:
21                    interaction = torch.zeros(self.embedding_dim, device
22                    =tag_vectors.device)
23                else:
24                    tag_embed = self.tag_embeddings(active_tags) # (
25                    num_active_tags, embedding_dim)
26                    # Compute outer product (sum of pairwise
27                    interactions)
28                    interaction = torch.matmul(tag_embed.t(), tag_embed)
29                    # (embedding_dim, embedding_dim)
30                    interaction = interaction.sum(dim=1) # Sum over
31                    interactions
32                    interaction_embeddings.append(interaction)
33            interaction_embeddings = torch.stack(interaction_embeddings)
34            # (batch_size, embedding_dim)
35            # Concatenate
36            embedded = torch.cat([embedded, interaction_embeddings], dim
37            =1) # (batch_size, 2*embedding_dim)
38

```

Listing 5: Corrected Interaction Embeddings Implementation

#### 4.4.4 Advantages and Considerations

##### Advantages:

- **Captures Interdependencies:** Explicitly models how pairs of tags interact, potentially uncovering complex relationships.
- **Improved Performance:** Can enhance model performance by leveraging tag dependencies.

##### Considerations:

- **Scalability:** Pairwise (or higher-order) interactions can lead to a combinatorial explosion of terms, especially with many tags.
- **Increased Dimensionality:** Concatenating interaction terms significantly increases the embedding size.
- **Implementation Complexity:** Requires careful handling to compute interactions efficiently, especially for large batches or many tags.
- **Potential Overfitting:** More parameters and higher-dimensional embeddings can lead to overfitting, especially with limited data.

### 4.5 Neural Network-Based Embeddings (MLP Approach)

#### 4.5.1 Conceptual Overview

The **MLP Embeddings** approach employs a small neural network (typically a two-layer MLP) to transform the binary tag vectors into continuous embeddings. This allows capturing non-linear relationships between tags, offering more flexibility than linear methods.

#### 4.5.2 Detailed Working

- **MLP Structure:** A sequence of linear layers with activation functions (e.g., ReLU) that learn to map the input tag vectors to embeddings.
- **Non-Linear Transformation:** By incorporating activation functions, the MLP can model complex, non-linear interactions between tags.
- **Mathematical Formulation:**

$$\text{TagEmbed}(f_i) = \text{MLP}(T_i)$$

Where MLP consists of multiple layers with non-linear activations.

### 4.5.3 Code Explanation

```
1 class MLPEmbedding(nn.Module):
2     def __init__(self, num_tags, embedding_dim, hidden_dim=32):
3         super(MLPEmbedding, self).__init__()
4         self.mlp = nn.Sequential(
5             nn.Linear(num_tags, hidden_dim),
6             nn.ReLU(),
7             nn.Linear(hidden_dim, embedding_dim)
8         )
9
10    def forward(self, tag_vectors):
11        embedded = self.mlp(tag_vectors.float()) # (batch_size,
12        embedding_dim)
13        return embedded
```

Listing 6: MLP Embeddings Implementation

- **Initialization (`__init__`):**

- `hidden_dim`: Number of neurons in the hidden layer (default is 32).
- `self.mlp`: A sequential model comprising:
  1. `nn.Linear(num_tags, hidden_dim)`: First linear layer mapping input tags to hidden dimension.
  2. `nn.ReLU()`: Non-linear activation function.
  3. `nn.Linear(hidden_dim, embedding_dim)`: Second linear layer mapping hidden dimension to final embedding.

- **Forward Pass (`forward`):**

- `tag_vectors`: Binary matrix (`batch_size, num_tags`).
- `tag_vectors.float()`: Converts binary tensor to floating-point.
- `self.mlp(tag_vectors.float())`: Passes through the MLP to obtain embeddings (`batch_size, embedding_dim`).

### 4.5.4 Advantages and Considerations

#### Advantages:

- **Non-Linear Relationships**: Capable of modeling complex, non-linear interactions between tags.
- **Flexibility**: The architecture can be adjusted (e.g., number of layers, hidden dimensions) to suit specific needs.
- **Expressiveness**: Can capture intricate patterns in the tag data, potentially leading to better performance.

#### Considerations:

- **Increased Parameters**: More parameters than linear methods, which can be problematic with large tag sets or limited data.

- **Risk of Overfitting:** Higher model complexity necessitates regularization techniques to prevent overfitting.
- **Computational Cost:** More computationally intensive due to multiple layers and non-linear operations.

## 4.6 Bilinear Embeddings

### 4.6.1 Conceptual Overview

**Bilinear Embeddings** use bilinear transformations to capture multiplicative interactions between tags and embedding dimensions. This method enhances the expressiveness of embeddings by allowing each tag to influence each dimension of the embedding space uniquely.

### 4.6.2 Detailed Working

- **Bilinear Transformation:** Each tag has its own transformation matrix that interacts with the embedding dimensions. The final embedding is a combination of these transformations.
- **Mathematical Formulation:**

$$\text{TagEmbed}(f_i) = T_i \cdot W$$

Where:

- $T_i$  is the binary tag vector for feature  $f_i$ .
- $W$  is a bilinear weight matrix of shape  $(\text{num\_tags}, \text{embedding\_dim})$ .

### 4.6.3 Code Explanation

```

1 class BilinearEmbedding(nn.Module):
2     def __init__(self, num_tags, embedding_dim):
3         super(BilinearEmbedding, self).__init__()
4         self.W = nn.Parameter(torch.randn(num_tags, embedding_dim))
5
6     def forward(self, tag_vectors):
7         # Tag vectors: (batch_size, num_tags)
8         # W: (num_tags, embedding_dim)
9         embedded = tag_vectors.float() @ self.W # (batch_size,
embedding_dim)
10        return embedded

```

Listing 7: Bilinear Embeddings Implementation

- **Initialization (`__init__`):**
  - `self.W`: A learnable bilinear weight matrix initialized randomly. Shape:  $(\text{num\_tags}, \text{embedding\_dim})$ .
- **Forward Pass (`forward`):**
  - `tag_vectors`: Binary matrix  $(\text{batch\_size}, \text{num\_tags})$ .

- `tag_vectors.float()`: Converts to floating-point.
- `embedded = tag_vectors.float() @ self.W`: Performs matrix multiplication to obtain embeddings (`batch_size`, `embedding_dim`).

**Note:** This implementation is similar to the **Linear Mapping** method but framed as a bilinear transformation, emphasizing the multiplicative interaction between tags and embedding dimensions.

#### 4.6.4 Advantages and Considerations

**Advantages:**

- **Interaction Modeling:** Captures multiplicative interactions between tags and embedding dimensions, enhancing expressiveness.
- **Expressive Power:** More capable of representing complex relationships compared to additive or simple linear methods.

**Considerations:**

- **Parameter Count:** Requires  $\text{num\_tags} \times \text{embedding\_dim}$  parameters, similar to linear mapping.
- **Training Complexity:** Bilinear transformations can be harder to train and may require careful initialization and optimization.
- **Overfitting Risk:** More parameters can lead to overfitting, especially with limited data.

### 4.7 Attention-Based Embeddings

#### 4.7.1 Conceptual Overview

**Attention-Based Embeddings** incorporate an attention mechanism to dynamically weigh the importance of each tag for a given feature. This allows the model to focus more on relevant tags while diminishing the influence of less important ones.

#### 4.7.2 Detailed Working

- **Attention Mechanism:** Computes attention weights for each tag based on the feature's context, then uses these weights to create a weighted sum of tag embeddings.
- **Dynamic Weighting:** Unlike static methods, attention-based embeddings adjust the importance of each tag during training, allowing for more nuanced representations.
- **Mathematical Formulation:**

$$\text{AttentionWeights} = \text{Softmax}(\text{Linear}(T_i \cdot E))$$

$$\text{TagEmbed}(f_i) = \sum_{j=1}^K \text{AttentionWeights}_j \cdot E_j$$

Where:



- $E_j$  is the embedding for tag  $j$ .
- $\text{AttentionWeights}_j$  is the attention weight for tag  $j$ .

### 4.7.3 Code Explanation

```

1 class AttentionEmbedding(nn.Module):
2     def __init__(self, num_tags, embedding_dim, attention_dim=32):
3         super(AttentionEmbedding, self).__init__()
4         self.tag_embeddings = nn.Embedding(num_tags, embedding_dim)
5         self.attention = nn.Sequential(
6             nn.Linear(embedding_dim, attention_dim),
7             nn.Tanh(),
8             nn.Linear(attention_dim, 1),
9             nn.Softmax(dim=1)
10        )
11
12    def forward(self, tag_vectors):
13        # tag_vectors: (batch_size, num_tags)
14        embeddings = self.tag_embeddings.weight # (num_tags,
embedding_dim)
15        embedded = tag_vectors @ embeddings # (batch_size,
embedding_dim)
16
17        # Calculate attention weights
18        attention_weights = self.attention(embedded) # (batch_size, 1)
19        # Weighted sum (here it's scalar attention, can be extended)
20        weighted = embedded * attention_weights # (batch_size,
embedding_dim)
21        return weighted

```

Listing 8: Attention-Based Embeddings Implementation

- **Initialization (`__init__`):**

- `attention_dim`: Dimension of the intermediate attention layer (default is 32).
- `self.tag_embeddings`: Embedding layer (`num_tags`, `embedding_dim`).
- `self.attention`: Sequential model comprising:
  1. `nn.Linear(embedding_dim, attention_dim)`: Projects embeddings to attention dimension.
  2. `nn.Tanh()`: Activation function.
  3. `nn.Linear(attention_dim, 1)`: Projects to a single attention score per tag.
  4. `nn.Softmax(dim=1)`: Normalizes attention scores across tags.

- **Forward Pass (`forward`):**

- `embeddings = self.tag_embeddings.weight`:  
Retrieves embedding matrix (`num_tags`, `embedding_dim`).
- `embedded = tag_vectors @ embeddings`:  
Computes summed embeddings (`batch_size`, `embedding_dim`).
- **Attention Weights:**

- \* `self.attention(embedded)`: Passes the summed embeddings through the attention network to obtain attention weights (`batch_size, 1`).
- \* `weighted = embedded * attention_weights`: Applies the attention weights to the embeddings. Since `attention_weights` is (`batch_size, 1`), broadcasting occurs across the embedding dimension.

**Note:** The current implementation applies a single scalar attention weight to the entire embedding vector. To implement per-tag attention, the code needs to be adjusted to compute attention weights for each tag individually.

#### 4.7.4 Correct Implementation for Per-Tag Attention

```

1 class AttentionEmbedding(nn.Module):
2     def __init__(self, num_tags, embedding_dim, attention_dim=32):
3         super(AttentionEmbedding, self).__init__()
4         self.tag_embeddings = nn.Embedding(num_tags, embedding_dim)
5         self.attention = nn.Sequential(
6             nn.Linear(embedding_dim, attention_dim),
7             nn.Tanh(),
8             nn.Linear(attention_dim, 1)
9         )
10
11     def forward(self, tag_vectors):
12         batch_size = tag_vectors.size(0)
13         num_tags = tag_vectors.size(1)
14         embeddings = self.tag_embeddings.weight # (num_tags,
15         embedding_dim)
16         embeddings = embeddings.unsqueeze(0).expand(batch_size, num_tags
17         , -1) # (batch_size, num_tags, embedding_dim)
18         tag_vectors = tag_vectors.unsqueeze(2).float() # (batch_size,
19         num_tags, 1)
20         # Mask inactive tags
21         masked_embeddings = embeddings * tag_vectors # (batch_size,
22         num_tags, embedding_dim)
23         # Compute attention scores for each tag
24         attention_scores = self.attention(masked_embeddings.view(-1,
25         embeddings.size(-1))) # (batch_size*num_tags, 1)
26         attention_scores = attention_scores.view(batch_size, num_tags)
27         # (batch_size, num_tags)
28         # Apply mask to zero out inactive tags
29         attention_scores = attention_scores * tag_vectors.squeeze(2) #
30         (batch_size, num_tags)
31         # Re-normalize attention scores
32         attention_weights = nn.functional.softmax(attention_scores, dim
33         =1) # (batch_size, num_tags)
34         attention_weights = attention_weights.unsqueeze(2) # (
35         batch_size, num_tags, 1)
36         # Weighted sum of embeddings
37         weighted = (masked_embeddings * attention_weights).sum(dim=1) #
38         (batch_size, embedding_dim)
39         return weighted

```

Listing 9: Corrected Attention-Based Embeddings Implementation

### 4.7.5 Advantages and Considerations

#### Advantages:

- **Dynamic Weighting:** Adjusts the influence of each tag based on the feature's context.
- **Focus on Relevant Tags:** Enhances the embedding by emphasizing more important tags.
- **Improved Performance:** Potentially leads to better representations by prioritizing significant tags.

#### Considerations:

- **Increased Complexity:** More layers and operations, making the model more complex.
- **Potential Training Instability:** Attention mechanisms can be sensitive to hyperparameters and may require careful tuning.
- **Computational Overhead:** Additional computations for calculating attention weights.

## 4.8 Hierarchical Embeddings

### 4.8.1 Conceptual Overview

**Hierarchical Embeddings** are designed for scenarios where tags have an inherent hierarchical structure (e.g., categories and subcategories). This method leverages the hierarchy to inform the embedding process, allowing the model to capture relationships at different levels of granularity.

### 4.8.2 Detailed Working

- **Hierarchy Structure:** Tags are organized into parent-child relationships, where parent tags encompass one or more child tags.
- **Embedding Incorporation:** The embeddings of child tags are combined with their parent tags' embeddings, allowing the model to inherit hierarchical information.
- **Mathematical Formulation:**

$$E'_c = E_c + E_p$$

Where  $E_c$  is the original embedding of the child tag and  $E_p$  is the embedding of the parent tag.

### 4.8.3 Code Explanation

```
1 class HierarchicalEmbedding(nn.Module):
2     def __init__(self, num_tags, embedding_dim, hierarchy):
3         """
4         Args:
5             num_tags (int): Total number of unique tags.
6             embedding_dim (int): Dimension of the embedding vectors.
7             hierarchy (dict): Dictionary defining parent-child
8             relationships among tags.
9             e.g., {parent_tag_idx: [child_tag_idx1,
10             child_tag_idx2, ...], ...}
11         """
12         super(HierarchicalEmbedding, self).__init__()
13         self.tag_embeddings = nn.Embedding(num_tags, embedding_dim)
14         self.hierarchy = hierarchy # e.g., {parent_tag: [child_tags]}
15
16     def forward(self, tag_vectors):
17         batch_size = tag_vectors.size(0)
18         embeddings = self.tag_embeddings.weight # (num_tags,
19         embedding_dim)
20         combined = tag_vectors @ embeddings # (batch_size,
21         embedding_dim)
22
23         # Incorporate hierarchical information
24         for parent, children in self.hierarchy.items():
25             parent_embedding = self.tag_embeddings(torch.tensor(parent,
26             device=tag_vectors.device)) # (embedding_dim,)
27             children_tensor = torch.tensor(children, device=tag_vectors.
28             device) # (num_children,)
29             children_active = tag_vectors[:, children].float() # (
30             batch_size, num_children)
31             children_sum = children_active.sum(dim=1) # (
32             batch_size,)
33             combined += children_sum.unsqueeze(1) * parent_embedding #
34             (batch_size, embedding_dim)
35         return combined
```

Listing 10: Hierarchical Embeddings Implementation

- **Initialization (`__init__`):**

- `hierarchy`: A dictionary mapping parent tag indices to lists of child tag indices.
- `self.tag_embeddings`: Embedding layer (`num_tags`, `embedding_dim`).

- **Forward Pass (`forward`):**

- `tag_vectors`: Binary matrix (`batch_size`, `num_tags`).
- `embeddings = self.tag_embeddings.weight`:  
Retrieves embedding matrix.
- `combined = tag_vectors @ embeddings`:  
Computes summed embeddings (`batch_size`, `embedding_dim`).
- **Hierarchical Information Incorporation:**

- \* Iterates over each parent and its children.
- \* `parent_embedding`: Retrieves the embedding for the parent tag.
- \* `children_tensor`: Converts child tag indices to a tensor.
- \* `children_active`: Extracts the active child tags for each sample (`batch_size`, `num_children`).
- \* `children_sum`: Sums the active child tags for each sample (`batch_size`,).
- \* `combined += children_sum.unsqueeze(1) * parent_embedding`: Adds the parent embedding scaled by the number of active child tags to the combined embeddings.

**Note:** This implementation assumes that for each active child tag, the parent tag’s embedding should be added once. Depending on the specific hierarchy and desired behavior, alternative methods of combining parent and child embeddings might be more appropriate.

#### 4.8.4 Advantages and Considerations

##### Advantages:

- **Captures Structured Relationships:** Leverages hierarchical structures to inform embeddings, enhancing representational power.
- **Improved Generalization:** Hierarchical information can help in better generalizing to unseen tag combinations.
- **Hierarchical Consistency:** Ensures that child tags inherit information from their parent categories.

##### Considerations:

- **Requires Hierarchical Structure:** Only applicable if tags have a meaningful hierarchy.
- **Implementation Complexity:** More involved to implement, especially with deep or complex hierarchies.
- **Potential Redundancy:** If not managed carefully, adding parent embeddings might lead to redundant information.

## 4.9 Pretrained Embeddings

### 4.9.1 9.1 Conceptual Overview

**Pretrained Embeddings** leverage embeddings learned from external sources or larger datasets. If such embeddings exist for your tags, they can be incorporated into your model to provide a head start, especially beneficial when training data is limited.

### 4.9.2 Detailed Working

- **Pretrained Embedding Matrix:** A precomputed embedding matrix, often obtained from models trained on similar tasks or large-scale data.
- **Integration:** The pretrained embeddings are loaded into the embedding layer, optionally allowing them to be fine-tuned during model training.
- **Mathematical Formulation:**

$$\text{TagEmbed}(f_i) = \text{PretrainedE}_j$$

Where  $\text{PretrainedE}_j$  is the pretrained embedding for tag  $j$ .

### 4.9.3 Code Explanation

```
1 class PretrainedEmbedding(nn.Module):
2     def __init__(self, pretrained_embeddings):
3         """
4         Args:
5             pretrained_embeddings (torch.Tensor): Pretrained embedding
6             matrix of shape (num_tags, embedding_dim).
7         """
8         super(PretrainedEmbedding, self).__init__()
9         self.tag_embeddings = nn.Embedding.from_pretrained(
10             pretrained_embeddings, freeze=False)
11
12     def forward(self, tag_vectors):
13         # tag_vectors: (batch_size, num_tags)
14         # W_tags: (num_tags, embedding_dim)
15         embedded = tag_vectors.float() @ self.tag_embeddings.weight # (
16             batch_size, embedding_dim)
17         return embedded
```

Listing 11: Pretrained Embeddings Implementation

- **Initialization (`__init__`):**
  - `pretrained_embeddings`: A tensor containing pretrained embeddings of shape `(num_tags, embedding_dim)`.
  - `self.tag_embeddings = nn.Embedding.from_pretrained(...)`: Initializes the embedding layer with the pretrained embeddings. `freeze=False` allows these embeddings to be fine-tuned during training.
- **Forward Pass (`forward`):**
  - `tag_vectors`: Binary matrix `(batch_size, num_tags)`.
  - `self.tag_embeddings.weight`: Retrieves the pretrained embedding matrix.
  - `embedded = tag_vectors.float() @ self.tag_embeddings.weight`: Multiplies binary tag vectors with the embedding matrix to obtain summed embeddings `(batch_size, embedding_dim)`.

#### 4.9.4 Advantages and Considerations

##### Advantages:

- **Leverage External Knowledge:** Utilizes embeddings trained on larger or domain-specific datasets, potentially enhancing performance.
- **Faster Convergence:** Provides a good initialization, leading to faster training and better generalization.
- **Improved Representations:** Pretrained embeddings often capture semantic relationships that can benefit downstream tasks.

##### Considerations:

- **Availability:** Pretrained embeddings might not be available for all types of tags or specific domains.
- **Compatibility:** Ensure that the dimensions and semantics of pretrained embeddings align with your model's requirements.
- **Fine-Tuning:** Deciding whether to freeze the embeddings or allow them to be fine-tuned is crucial and depends on the size of your dataset and similarity to the pretrained data.

## 5 Choosing the Right Embedding Method

Selecting the appropriate embedding method depends on several factors, including the nature of your tags, the complexity of their relationships, computational resources, and the specific requirements of your model. Here's a guide to help you decide:

### 1. Simplicity vs. Expressiveness:

- **Simple Methods:** Start with **Sum of Embeddings** or **Linear Mapping** for their simplicity and efficiency.
- **Complex Relationships:** If tags have intricate interactions, consider **Interaction Embeddings** or **Attention-Based Embeddings**.

### 2. Data Availability:

- **Limited Data:** Simpler methods like **Sum of Embeddings** are less prone to overfitting.
- **Abundant Data:** More complex methods can be explored to capture deeper relationships.

### 3. Hierarchy and Structure:

- If your tags are hierarchical, **Hierarchical Embeddings** can capture structured relationships effectively.

### 4. Pretrained Resources:

- Utilize **Pretrained Embeddings** if relevant embeddings are available to incorporate external knowledge.

## 5. Computational Resources:

- Consider the computational cost of each method. Simpler methods are generally faster and require less memory.

## 6. Hybrid Approaches:

- Combine multiple embedding techniques, such as using both **Sum of Embeddings** and **Attention-Based Embeddings**, to leverage their respective strengths.

# 6 Integrating Embeddings into the TimeXer Model

To effectively incorporate feature group embeddings into the TimeXer model, follow these steps:

1. **Embed Feature Tags:** Use the selected embedding method to convert binary tag vectors into continuous embeddings.
2. **Combine with Feature Data:** Integrate the tag embeddings with the original feature representations. This can be done via concatenation or addition.
3. **Modify the TimeXer Model:** Adjust the input dimensions of the TimeXer model to accommodate the additional embedding dimensions.

## 6.1 Example Integration Using Sum of Embeddings

```

1 import torch
2 import torch.nn as nn
3
4 class TimeXerWithSumEmbedding(nn.Module):
5     def __init__(self, input_dim, projected_dim, hidden_dim, num_layers,
6         output_dim, num_tags, embedding_dim):
7         super(TimeXerWithSumEmbedding, self).__init__()
8         self.tag_embedding = SumEmbedding(num_tags, embedding_dim)
9         # Adjust input dimension to include embedding dimension
10        self.projection = nn.Linear(input_dim + embedding_dim,
11            projected_dim)
12        self.dropout = nn.Dropout(p=0.1)
13        self.transformer_layer = nn.TransformerEncoder(
14            nn.TransformerEncoderLayer(d_model=projected_dim, nhead=4,
15                dim_feedforward=hidden_dim, batch_first=True),
16                num_layers=num_layers
17        )
18        self.fc = nn.Linear(projected_dim, output_dim)
19
20    def forward(self, x, tag_vectors):
21        # Embed tags
22        embedded_tags = self.tag_embedding(tag_vectors) # (batch_size,
23            embedding_dim)
24        # Expand embedded_tags to match sequence length

```



```

21     embedded_tags = embedded_tags.unsqueeze(1).repeat(1, x.size(1),
1) # (batch_size, seq_len, embedding_dim)
22     # Concatenate with input features
23     x = torch.cat([x, embedded_tags], dim=-1) # (batch_size,
seq_len, input_dim + embedding_dim)
24     # Proceed with TimeXer layers
25     x = self.projection(x) # (batch_size, seq_len, projected_dim)
26     x = self.dropout(x)
27     x = self.transformer_layer(x)
28     x = self.fc(x)
29     return x

```

Listing 12: TimeXer Integration with Sum of Embeddings

### 6.1.1 Explanation

#### 1. Embedding Layer Initialization:

- `self.tag_embedding = SumEmbedding(num_tags, embedding_dim)`: Initializes the tag embedding layer using the **Sum of Embeddings** method.

#### 2. Adjusting Input Dimensions:

- `self.projection = nn.Linear(input_dim + embedding_dim, projected_dim)`: Concatenates the original feature dimensions with the embedding dimensions to align with the transformer's expected input.

#### 3. Forward Pass Steps:

1. **Embedding Tags**: Converts binary tag vectors into continuous embeddings.
2. **Broadcasting**: Expands the tag embeddings across the sequence length to match the temporal dimension.
3. **Concatenation**: Merges the original features with the tag embeddings, enriching the feature representation.
4. **Projection and Transformation**: Passes the combined features through the projection layer, dropout, transformer encoder, and finally the fully connected output layer.

### 6.1.2 Usage Example

```

1 # Define model parameters
2 input_dim = 10           # Number of original feature dimensions
3 projected_dim = 64       # Embedding size for transformer
4 hidden_dim = 128        # Hidden layer size in transformer
5 num_layers = 3           # Number of transformer layers
6 output_dim = 1          # Output dimension (e.g., regression)
7 num_tags = 17            # Number of unique tags
8 embedding_dim = 16       # Embedding dimension for tags
9
10 # Instantiate the model
11 model = TimeXerWithSumEmbedding(input_dim, projected_dim, hidden_dim,
num_layers, output_dim, num_tags, embedding_dim)
12

```

```

13 # Example input:
14 batch_size = 32
15 seq_length = 50
16 x = torch.randn(batch_size, seq_length, input_dim)           # Original
   features
17 tag_vectors = torch.randint(0, 2, (batch_size, num_tags)).float() #
   Binary tag vectors
18
19 # Forward pass
20 output = model(x, tag_vectors)
21
22 print(output.shape) # Expected: (32, 50, 1)

```

Listing 13: TimeXerWithSumEmbedding Usage Example

## 7 Final Recommendations and Best Practices

Choosing the right embedding method depends on your specific use case, dataset characteristics, and computational resources. Here are some guidelines to help you decide:

### 1. Start Simple:

- **Sum of Embeddings:** Given its simplicity and effectiveness, especially with a manageable number of tags (17), this method is a good starting point.
- **Linear Mapping:** If you require more expressiveness than summation offers, consider linear mapping.

### 2. Evaluate Model Performance:

- Compare the performance of different embedding methods using validation metrics.
- Monitor for overfitting, especially with more complex methods like **Linear Mapping** or **MLP Embeddings**.

### 3. Consider Alternative Methods if Needed:

- If initial experiments with **Sum of Embeddings** do not yield satisfactory performance, consider exploring **Linear Mapping** or **MLP Embeddings** as next steps, keeping in mind the trade-offs in complexity and overfitting risks.

### 4. Regularization Techniques:

- **Dropout:** Apply dropout to embedding layers and downstream layers to prevent overfitting.
- **L2 Regularization:** Implement L2 regularization on embedding weights to encourage smaller weights and reduce overfitting.

### 5. Hyperparameter Tuning:

- **Embedding Dimension ( $d_e$ ):** Experiment with different embedding dimensions (e.g., 16, 32) to find the optimal balance between expressiveness and computational efficiency.

- **Hidden Dimensions:** If using **MLP Embeddings**, tune the hidden dimensions to capture sufficient complexity without overfitting.

#### 6. **Visualization:**

- Use techniques like PCA or t-SNE to visualize learned embeddings, ensuring that they capture meaningful distinctions between tags despite their anonymity.

#### 7. **Scalability Considerations:**

- While 17 tags are manageable, anticipate potential increases in tag numbers. Ensure that your embedding method can scale accordingly without significant performance degradation.

#### 8. **Integration with Downstream Layers:**

- Ensure that the combined feature and embedding dimensions align with the expectations of downstream layers, such as transformers or other neural network components.

#### 9. **Experimentation:**

- Experiment with multiple embedding methods and evaluate their impact on model performance. Use cross-validation and ablation studies to guide your selection.

## 8 Conclusion

Feature Group Embeddings are pivotal in transforming categorical tag data into meaningful continuous representations that enhance machine learning models' performance. Among the various methods, **Sum of Embeddings** stands out for its simplicity, efficiency, and effectiveness, especially in scenarios with a manageable number of anonymized tags and features. By adhering to best practices and methodically evaluating each approach, you can optimize your models to leverage categorical information effectively, leading to more accurate and robust predictions.

## Appendix: Summary of Embedding Methods

Method	Description	Pros & Cons
Sum of Embeddings	Sum the embeddings of active tags	<b>Pros:</b> Simple, efficient, parameter-light <b>Cons:</b> May miss complex tag interactions
Linear Mapping	Learn a weight matrix to map tag vectors	<b>Pros:</b> Captures interactions, more expressive <b>Cons:</b> More parameters, risk of overfitting
Multi-Head Embeddings	Multiple embeddings per tag to capture different aspects	<b>Pros:</b> Enhanced representation, parallel learning <b>Cons:</b> Increased complexity and computation
Interaction Embeddings	Explicitly model interactions between tags	<b>Pros:</b> Captures dependencies, improved performance <b>Cons:</b> Scalability issues with many tags, complex implementation
MLP Embeddings	Use a neural network to non-linearly transform tag vectors	<b>Pros:</b> Captures non-linear relationships, flexible <b>Cons:</b> More parameters, higher risk of overfitting
Bilinear Embeddings	Use bilinear transformations to model multiplicative interactions	<b>Pros:</b> Captures multiplicative interactions <b>Cons:</b> Increased parameters and complexity
Attention-Based Embeddings	Apply attention mechanisms to weigh tag importance	<b>Pros:</b> Dynamic weighting, focuses on relevant tags <b>Cons:</b> Adds computational overhead, training complexity
Hierarchical Embeddings	Leverage hierarchical relationships among tags	<b>Pros:</b> Captures structured relationships <b>Cons:</b> Requires hierarchical tag structure, complex implementation
Pretrained Embeddings	Utilize pretrained embeddings if available	<b>Pros:</b> Leverages external knowledge, potential performance boost <b>Cons:</b> Limited applicability, requires compatible pretrained embeddings

Table 1: Summary of Embedding Methods

## References

1. [PyTorch Documentation](#)
2. [Understanding Embeddings](#)
3. [Transformers and Attention Mechanisms](#)