Embedding Matrix

Create a trainable embedding matrix to calculate word embeddings.

Chapter Goals:

- Learn how the embedding matrix is initialized during training
- Create a function for the forward run of embedding training

A. Variable creation

In TensorFlow, we use the tf.layers.dense function to create each layer of a fully-connected neural network. This is a function provided by TensorFlow that automatically handles all the setup required for the neural network's weights. So previously all we needed to do was initialize the weight variables (via tf.global_variables_initializer) prior to running the computation graph.

However, for our embedding model, the variable we need to create is a 2-D matrix that contains the embedding vectors for each vocabulary word ID. Therefore, we need to manually create this variable using the tf.get_variable function.

The function only takes in one required argument, which is the name of the variable. Of the possible keyword arguments, a couple important ones to know are shape and dtype, which allow us to specify the shape and type, respectively, of the created variable.

Below we demonstrate example usages of tf.get_variable:



In the example above, the first variable is created with just the shape
keyword argument (so its default type is tf.float32), while the second is
manually set to type tf.int64 (the _ref suffix does not affect the type).

During training, the first call to <code>tf.get_variable</code> with a given variable name will create a new variable. Each subsequent call to the function using the same variable name will retrieve the existing variable, rather than create a new one. This allows us to continue training with the same embedding matrix.

B. Variable initializers

Another useful keyword argument to know for tf.get_variable is the initializer argument, which is used to specify how the variable will be initialized. The argument can also be used to specify the variable's shape (in which case we don't need to set the shape keyword argument).

The default initializer for variables created by <code>tf.get_variable</code> is a <code>glorot_uniform_initializer</code>. However, we can use our own probability distribution to randomly initialize the variable. A popular probability distribution to use is <code>tf.random_uniform</code>, corresponding to a uniform distribution.

Below we demonstrate an example usage of tf.random_uniform to initialize a variable:

```
import tensorflow as tf
init = tf.random_uniform((5, 10),minval=-1,maxval=2)
v = tf.get_variable('v1', initializer=init)
```

In the example above, we used <code>tf.random_uniform</code> to return a tensor with shape <code>(5, 10)</code>, containing randomly chosen values from the uniform distribution in the range <code>[-1, 2]</code>. We then used <code>init</code> to initialize <code>v</code> with its randomly chosen values.

C. Embedding lookup

When training the embedding model, the "forward" run consists of variable initialization/retrieval followed by *embedding lookup* for the current

iteration's training batch. Embedding lookup refers to retrieving the

embedding vectors for each word ID in the training batch. Since the embedding matrix's rows are each unique embedding vectors, we perform the lookup simply by retrieving the rows corresponding to the training batch's word IDs.

The function we use to retrieve the embedding vectors is <code>tf.nn.embedding_lookup</code>. It takes in two required arguments, which are the embedding matrix variable and vocabulary IDs to lookup.

Below we demonstrate an example usage of tf.nn.embedding_lookup:

```
import tensorflow as tf
emb_mat = tf.get_variable('v1', shape=(5, 10))
word_ids = tf.constant([0, 3])
emb_vecs = tf.nn.embedding_lookup(emb_mat, word_ids)
print(emb_vecs)
```

In the example above, we used <code>tf.nn.embedding_lookup</code> to retrieve embedding vectors from the embedding matrix, <code>emb_mat</code>, for the word IDs <code>0</code> and <code>3</code>. The output tensor from the embedding lookup will contain the embedding vector for ID <code>0</code> in the first row, and the embedding vector for ID <code>3</code> in the second row.

Time to Code!

In this chapter, you'll be completing the **forward** function, which retrieves embeddings using an embedding matrix.

First, you'll complete the helper function <code>get_initializer</code>. This function returns the initializer for the embedding matrix.

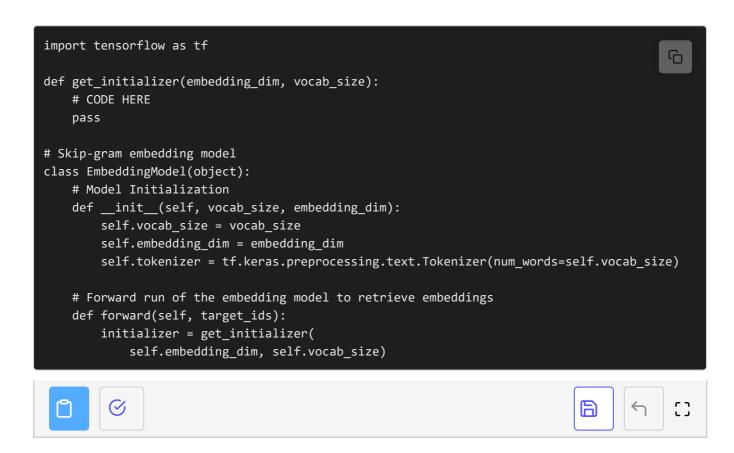
We're going to initialize our embedding matrix variable with a random uniform initializer. As a rule of thumb, we'll set the uniform distribution bounds to [-0.5/d, 0.5/d), where d is the embedding vector size (i.e. embedding dim).

Set initial_bounds equal to 0.5 divided by embedding_dim.

We want our embedding matrix to have shape (vocab_size, embedding_dim), since each row of the matrix corresponds to an embedding vector for a unique vocabulary word.

Set initializer equal to tf.random_uniform with the specified matrix shape as the required argument.

The minval and maxval keyword arguments should correspond to negative and positive initial_bounds, respectively.



Now you'll work in the main function, forward (where it says "CODE HERE"). Using the initializer, we can create/retrieve the embedding matrix variable, then use it to get embeddings.

Set self.embedding_matrix equal to tf.get_variable with 'embedding_matrix' as the required argument and initializer as the initializer keyword argument.

Set embeddings equal to tf.nn.embedding_lookup with self.embedding_matrix and target_ids as the two required arguments. Then return embeddings.

```
# Skip-gram embedding model
class EmbeddingModel(object):
    # Model Initialization

def __init__(self, vocab_size, embedding_dim):
    self.vocab_size = vocab_size
    self.embedding_dim = embedding_dim
    self.tokenizer = tf.keras.preprocessing.text.Tokenizer(num_words=self.vocab_size)

# Forward run of the embedding model to retrieve embeddings
def forward(self, target_ids):
    initializer = get_initializer(
        self.embedding_dim, self.vocab_size)
    # CODE HERE
```









[]