

-0.75

Figure 12. A one-hot encoding of hot dog provided as input to a deep neural network. An embedding layer translates the one-hot encoding into the three-dimensional embedding vector [2.98, -0.75, 0].

closer to each other. As previously mentioned, the dimensions that an actual model chooses for its embeddings are unlikely to be as intuitive or understandable as in this example.

In the course of training, the weights of the embedding layer will be optimized so that the embedding vectors for similar examples are

Contextual embeddings

One limitation of word2vec static embedding vectors is that words can mean different things in different contexts. "Yeah" means one thing on its own, but the opposite in the phrase "Yeah, right." "Post" can mean "mail," "to put in the mail," "earring backing," "marker at the end of a horse race," "postproduction," "pillar," "to put up a notice," "to station a guard or soldier," or "after," among other possibilities.

However, with static embeddings, each word is represented by a single point in vector space, even though it may have a variety of meanings.

In the last exercise, you discovered the limitations of static embeddings for the word orange, which can signify either a color or a type of

fruit. With only one static embedding, orange will always be closer to other colors than to juice when trained on the word2vec dataset. Contextual embeddings were developed to address this limitation. Contextual embeddings allow a word to be represented by multiple embeddings that incorporate information about the surrounding words as well as the word itself. Orange would have a different embedding

Some methods for creating contextual embeddings, like ELMo, take the static embedding of an example, such as the word2vec vector for a word in a sentence, and transform it by a function that incorporates information about the words around it. This produces a contextual embedding.

• For ELMo models specifically, the static embedding is aggregated with embeddings taken from other layers, which encode front-to-

Click here for details on contextual embeddings

for every unique sentence containing the word in the dataset.

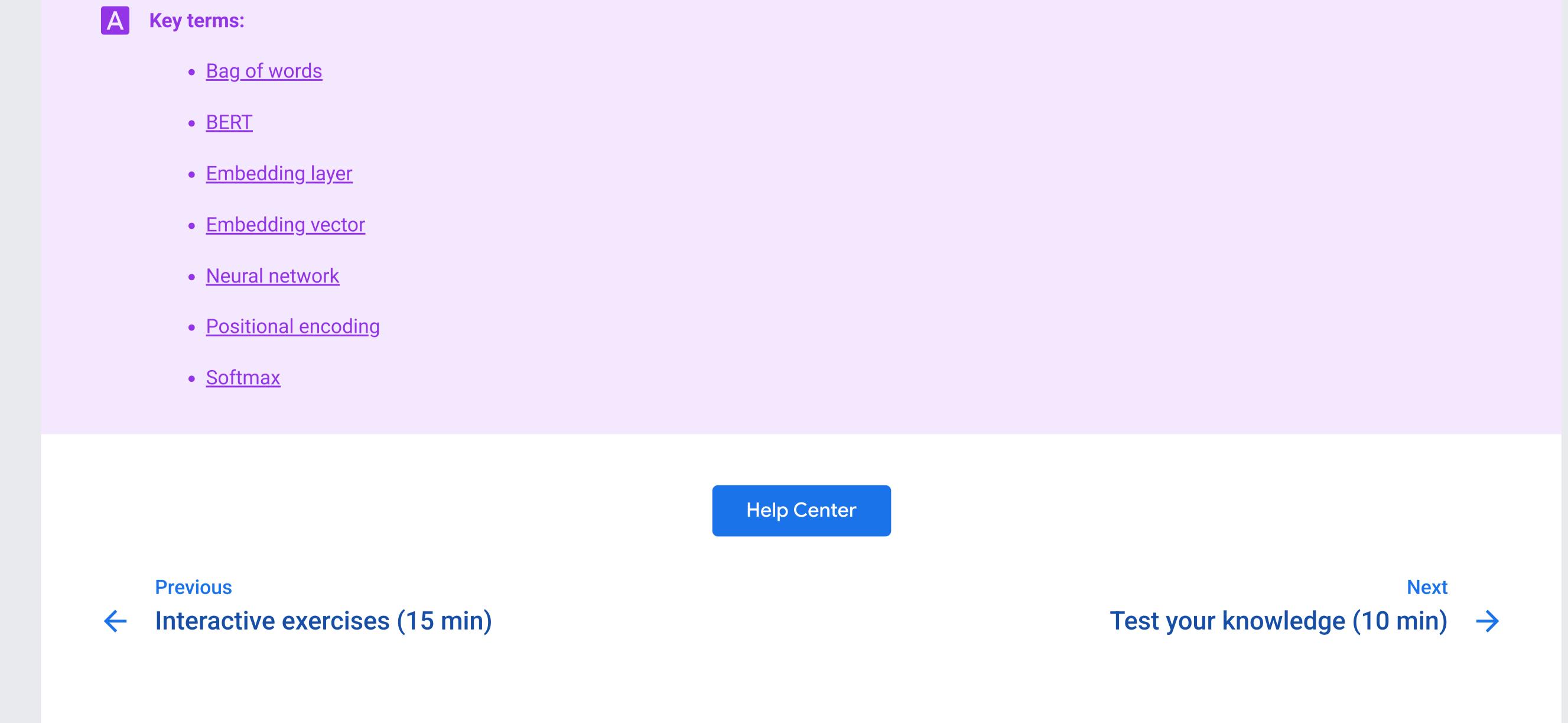
- back and back-to-front readings of the sentence. • BERT models mask part of the sequence that the model takes as input.
- Transformer models use a self-attention layer to weight the relevance of the other words in a sequence to each individual word. They

details, see the Google Developers Site Policies. Java is a registered trademark of Oracle and/or its affiliates.

also add the relevant column from a positional embedding matrix (see positional encoding) to each previously learned token embedding, element by element, to produce the input embedding that is fed into the rest of the model for inference. This **input** embedding, unique to each distinct textual sequence, is a contextual embedding.

```
NOTE: See the <u>LLM module</u> for more details on transformers and encoder-decoder architecture.
```

While the models described above are language models, contextual embeddings are useful in other generative tasks, like images. An embedding of the pixel RGB values in a photo of a horse provides more information to the model when combined with a positional matrix representing each pixel and some encoding of the neighboring pixels, creating contextual embeddings, than the original static embeddings of the RGB values alone.



Last updated 2025-04-10 UTC.

Except as otherwise noted, the content of this page is licensed under the Creative Commons Attribution 4.0 License, and code samples are licensed under the Apache 2.0 License. For

Was this helpful?

Send feedback

