

## tf.contrib.seq2seq.LuongAttention

## Contents

Class LuongAttention

## Properties

alignments\_size

batch\_size

Class **LuongAttention**

Defined in [tensorflow/contrib/seq2seq/python/ops/attention\\_wrapper.py](#).

See the guide: [Seq2seq Library \(contrib\) > Attention](#)

Implements Luong-style (multiplicative) attention scoring.

This attention has two forms. The first is standard Luong attention, as described in:

Minh-Thang Luong, Hieu Pham, Christopher D. Manning. "Effective Approaches to Attention-based Neural Machine Translation." EMNLP 2015. <https://arxiv.org/abs/1508.04025>

The second is the scaled form inspired partly by the normalized form of Bahdanau attention.

To enable the second form, construct the object with parameter `scale=True`.

## Properties

**alignments\_size****batch\_size****keys****memory\_layer****query\_layer****values**

## Methods

**\_\_init\_\_**

```

__init__(
    num_units,
    memory,
    memory_sequence_length=None,
    scale=False,
    probability_fn=None,
    score_mask_value=float('-inf'),
    name='LuongAttention'
)

```

Construct the AttentionMechanism mechanism.

Args:

- `num_units` : The depth of the attention mechanism.
- `memory` : The memory to query; usually the output of an RNN encoder. This tensor should be shaped `[batch_size, max_time, ...]`.
- `memory_sequence_length` : (optional) Sequence lengths for the batch entries in memory. If provided, the memory tensor rows are masked with zeros for values past the respective sequence lengths.
- `scale` : Python boolean. Whether to scale the energy term.
- `probability_fn` : (optional) A `callable`. Converts the score to probabilities. The default is `tf.nn.softmax`. Other options include `tf.contrib.seq2seq.hardmax` and `tf.contrib.sparsemax.sparsemax`. Its signature should be: `probabilities = probability_fn(score)`.
- `score_mask_value` : (optional) The mask value for score before passing into `probability_fn`. The default is -inf. Only used if `memory_sequence_length` is not None.
- `name` : Name to use when creating ops.

**`__call__`**

```

__call__(
    query,
    previous_alignments
)

```

Score the query based on the keys and values.

Args:

- `query` : Tensor of dtype matching `self.values` and shape `[batch_size, query_depth]`.
- `previous_alignments` : Tensor of dtype matching `self.values` and shape `[batch_size, alignments_size]` (`alignments_size` is memory's `max_time`).

Returns:

- `alignments` : Tensor of dtype matching `self.values` and shape `[batch_size, alignments_size]` (`alignments_size` is memory's `max_time`).

**`initial_alignments`**

```
initial_alignments(  
    batch_size,  
    dtype  
)
```

Creates the initial alignment values for the `AttentionWrapper` class.

This is important for AttentionMechanisms that use the previous alignment to calculate the alignment at the next time step (e.g. monotonic attention).

The default behavior is to return a tensor of all zeros.

Args:

- `batch_size`: `int32` scalar, the batch\_size.
- `dtype`: The `dtype`.

Returns:

A `dtype` tensor shaped `[batch_size, alignments_size]` (`alignments_size` is the values' `max_time`).

---

*Except as otherwise noted, the content of this page is licensed under the [Creative Commons Attribution 3.0 License](#), and code samples are licensed under the [Apache 2.0 License](#). For details, see our [Site Policies](#). Java is a registered trademark of Oracle and/or its affiliates.*

*Last updated November 2, 2017.*

## Stay Connected

[Blog](#)

[GitHub](#)

[Twitter](#)

## Support

[Issue Tracker](#)

[Release Notes](#)

[Stack Overflow](#)

**English**

[Terms](#) | [Privacy](#)