


5100e359ae ▾

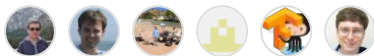
...

tensorflow / tensorflow / core / kernels / string_ngrams_op.cc

 jpinaar Update more int64 uses to int64_t ... ✖

 History

6 contributors



251 lines (225 sloc) | 10.1 KB

...

```

1  /* Copyright 2019 The TensorFlow Authors. All Rights Reserved.
2
3  Licensed under the Apache License, Version 2.0 (the "License");
4  you may not use this file except in compliance with the License.
5  You may obtain a copy of the License at
6
7      http://www.apache.org/licenses/LICENSE-2.0
8
9  Unless required by applicable law or agreed to in writing, software
10 distributed under the License is distributed on an "AS IS" BASIS,
11 WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
12 See the License for the specific language governing permissions and
13 limitations under the License.
14 =====*/
15
16 #include <locale>
17 #include <string>
18
19 #include "absl/strings/ascii.h"
20 #include "absl/strings/str_cat.h"
21 #include "tensorflow/core/framework/op_kernel.h"
22 #include "tensorflow/core/platform/errors.h"
23
24 namespace tensorflow {
25 namespace text {
26
27 namespace {
28 template <typename SPLITS_TYPE>
29 class StringNGramsOp : public tensorflow::OpKernel {

```

```

30 public:
31     explicit StringNGramsOp(tensorflow::OpKernelConstruction* context)
32         : tensorflow::OpKernel(context) {
33         OP_REQUIRES_OK(context, context->GetAttr("separator", &separator_));
34         OP_REQUIRES_OK(context, context->GetAttr("ngram_widths", &ngram_widths_));
35         OP_REQUIRES_OK(context, context->GetAttr("left_pad", &left_pad_));
36         OP_REQUIRES_OK(context, context->GetAttr("right_pad", &right_pad_));
37         OP_REQUIRES_OK(context, context->GetAttr("pad_width", &pad_width_));
38         OP_REQUIRES_OK(context, context->GetAttr("preserve_short_sequences",
39             &preserve_short_));
40     }
41
42     int get_pad_width(const int ngram_width) const {
43         // Ngrams can be padded with either a fixed pad width or a dynamic pad
44         // width depending on the 'pad_width' arg, but in no case should the padding
45         // ever be wider than 'ngram_width' - 1.
46         return std::min(pad_width_ < 0 ? ngram_width - 1 : pad_width_,
47             ngram_width - 1);
48     }
49
50     int get_num_ngrams(const int length, const int ngram_width) const {
51         int pad_width = get_pad_width(ngram_width);
52         return std::max(0, ((length + 2 * pad_width) - ngram_width) + 1);
53     }
54
55     void Compute(tensorflow::OpKernelContext* context) override {
56         for (int ngram_width : ngram_widths_) {
57             OP_REQUIRES(
58                 context, ngram_width > 0,
59                 errors::InvalidArgument("ngram_widths must contain positive values"));
60         }
61
62         const tensorflow::Tensor* data;
63         OP_REQUIRES_OK(context, context->input("data", &data));
64         const auto& input_data = data->flat<tstring>().data();
65
66         const tensorflow::Tensor* splits;
67         OP_REQUIRES_OK(context, context->input("data_splits", &splits));
68         const auto& splits_vec = splits->flat<SPLITS_TYPE>();
69
70         // Validate that the splits are valid indices into data, only if there are
71         // splits specified.
72         const int input_data_size = data->flat<tstring>().size();
73         const int splits_vec_size = splits_vec.size();
74         if (splits_vec_size > 0) {
75             int prev_split = splits_vec(0);
76             OP_REQUIRES(context, prev_split == 0,
77                 errors::InvalidArgument("First split value must be 0, got ",
78                     prev_split));

```

```

79     for (int i = 1; i < splits_vec_size; ++i) {
80         bool valid_splits = splits_vec(i) >= prev_split;
81         valid_splits = valid_splits && (splits_vec(i) <= input_data_size);
82         OP_REQUIRES(context, valid_splits,
83                     errors::InvalidArgument(
84                         "Invalid split value ", splits_vec(i), ", must be in [",
85                         prev_split, ", ", input_data_size, "]"));
86         prev_split = splits_vec(i);
87     }
88     OP_REQUIRES(context, prev_split == input_data_size,
89                 errors::InvalidArgument(
90                     "Last split value must be data size. Expected ",
91                     input_data_size, ", got ", prev_split));
92 }
93
94 int num_batch_items = splits_vec.size() - 1;
95 tensorflow::Tensor* ngrams_splits;
96 OP_REQUIRES_OK(
97     context, context->allocate_output(1, splits->shape(), &ngrams_splits));
98 auto ngrams_splits_data = ngrams_splits->flat<SPLITS_TYPE>().data();
99
100 // If there is no data or size, return an empty RT.
101 if (data->flat<tstring>().size() == 0 || splits_vec.size() == 0) {
102     tensorflow::Tensor* empty;
103     OP_REQUIRES_OK(context,
104                     context->allocate_output(0, data->shape(), &empty));
105     for (int i = 0; i <= num_batch_items; ++i) {
106         ngrams_splits_data[i] = 0;
107     }
108     return;
109 }
110
111 ngrams_splits_data[0] = 0;
112 for (int i = 1; i <= num_batch_items; ++i) {
113     int length = splits_vec(i) - splits_vec(i - 1);
114     int num_ngrams = 0;
115     for (int ngram_width : ngram_widths_)
116         num_ngrams += get_num_ngrams(length, ngram_width);
117     if (preserve_short_ && length > 0 && num_ngrams == 0) {
118         num_ngrams = 1;
119     }
120     ngrams_splits_data[i] = ngrams_splits_data[i - 1] + num_ngrams;
121 }
122
123 tensorflow::Tensor* ngrams;
124 OP_REQUIRES_OK(
125     context,
126     context->allocate_output(
127         0, TensorShape({ngrams_splits_data[num_batch_items]}), &ngrams));

```

```

128     auto ngrams_data = ngrams->flat<tstring>().data();
129
130     for (int i = 0; i < num_batch_items; ++i) {
131         auto data_start = &input_data[splits_vec(i)];
132         int output_start_idx = ngrams_splits_data[i];
133         for (int ngram_width : ngram_widths_) {
134             auto output_start = &ngrams_data[output_start_idx];
135             int length = splits_vec(i + 1) - splits_vec(i);
136             int num_ngrams = get_num_ngrams(length, ngram_width);
137             CreateNgrams(data_start, output_start, num_ngrams, ngram_width);
138             output_start_idx += num_ngrams;
139         }
140         // If we're preserving short sequences, check to see if no sequence was
141         // generated by comparing the current output start idx to the original
142         // one (ngram_splits_data). If no ngrams were generated, then they will
143         // be equal (since we increment output_start_idx by num_ngrams every
144         // time we create a set of ngrams.)
145         if (preserve_short_ && output_start_idx == ngrams_splits_data[i]) {
146             int data_length = splits_vec(i + 1) - splits_vec(i);
147             // One legitimate reason to not have any ngrams when preserve_short_
148             // is true is if the sequence itself is empty. In that case, move on.
149             if (data_length == 0) {
150                 continue;
151             }
152             // We don't have to worry about dynamic padding sizes here: if padding
153             // was dynamic, every sequence would have had sufficient padding to
154             // generate at least one ngram.
155             int ngram_width = data_length + 2 * pad_width_;
156             auto output_start = &ngrams_data[output_start_idx];
157             int num_ngrams = 1;
158             CreateNgrams(data_start, output_start, num_ngrams, ngram_width);
159         }
160     }
161 }
162
163 void CreateNgrams(const tstring* data, tstring* output, int num_ngrams,
164                 int ngram_width) const {
165     for (int ngram_index = 0; ngram_index < num_ngrams; ++ngram_index) {
166         int pad_width = get_pad_width(ngram_width);
167         int left_padding = std::max(0, pad_width - ngram_index);
168         int right_padding =
169             std::max(0, pad_width - (num_ngrams - (ngram_index + 1)));
170         int num_tokens = ngram_width - (left_padding + right_padding);
171         int data_start_index = left_padding > 0 ? 0 : ngram_index - pad_width;
172
173         // Calculate the total expected size of the ngram so we can reserve the
174         // correct amount of space in the string.
175         int ngram_size = 0;
176         // Size of the left padding.

```

```

177     ngram_size += left_padding * left_pad_.length();
178     // Size of the tokens.
179     for (int n = 0; n < num_tokens; ++n) {
180         ngram_size += data[data_start_index + n].length();
181     }
182     // Size of the right padding.
183     ngram_size += right_padding * right_pad_.length();
184     // Size of the separators.
185     int num_separators = left_padding + right_padding + num_tokens - 1;
186     ngram_size += num_separators * separator_.length();
187
188     // Build the ngram.
189     tstring* ngram = &output[ngram_index];
190     ngram->reserve(ngram_size);
191     for (int n = 0; n < left_padding; ++n) {
192         ngram->append(left_pad_);
193         ngram->append(separator_);
194     }
195     // Only output first num_tokens - 1 pairs of data and separator
196     for (int n = 0; n < num_tokens - 1; ++n) {
197         ngram->append(data[data_start_index + n]);
198         ngram->append(separator_);
199     }
200     // Handle case when there are no tokens or no right padding as these can
201     // result in consecutive separators.
202     if (num_tokens > 0) {
203         // If we have tokens, then output last and then pair each separator with
204         // the right padding that follows, to ensure ngram ends either with the
205         // token or with the right pad.
206         ngram->append(data[data_start_index + num_tokens - 1]);
207         for (int n = 0; n < right_padding; ++n) {
208             ngram->append(separator_);
209             ngram->append(right_pad_);
210         }
211     } else {
212         // If we don't have tokens, then the last item inserted into the ngram
213         // has been the separator from the left padding loop above. Hence,
214         // output right pad and separator and make sure to finish with a
215         // padding, not a separator.
216         for (int n = 0; n < right_padding - 1; ++n) {
217             ngram->append(right_pad_);
218             ngram->append(separator_);
219         }
220         ngram->append(right_pad_);
221     }
222
223     // In debug mode only: validate that we've reserved enough space for the
224     // ngram.
225     DCHECK_EQ(ngram_size, ngram->size());

```

```
226     }
227 }
228
229 string separator_;
230 string left_pad_;
231 string right_pad_;
232 bool use_pad_;
233 bool extend_pad_;
234 bool preserve_short_;
235
236 std::vector<int> ngram_widths_;
237 int pad_width_;
238 };
239
240 } // namespace
241 REGISTER_KERNEL_BUILDER(Name("StringNGrams")
242     .Device(tensorflow::DEVICE_CPU)
243     .TypeConstraint<int32>("Tsplits"),
244     StringNGramsOp<int32>);
245 REGISTER_KERNEL_BUILDER(Name("StringNGrams")
246     .Device(tensorflow::DEVICE_CPU)
247     .TypeConstraint<int64_t>("Tsplits"),
248     StringNGramsOp<int64_t>);
249
250 } // namespace text
251 } // namespace tensorflow
```