# 支持多值带权重、稀疏、共享embedding权重的DSSM召回实现（tensorflow2）

原创　xulu1352　数据与智能　今天

收录于话题

#推荐算法 3　#推荐系统 3　#TensorFlow 1　#深度学习 2　#召回 1

点击上方**"数据与智能"**，"星标或置顶公众号"

第一时间获取好内容



作者 | xulu1352　目前在一家互联网公司从事推荐算法工作

编辑 | lily

<div align="center">

| 0 |
|---|

</div>

<div align="center">

| 前序 |
|---|

</div>

　　关于DSSM模型原理及实现，网上已经有很多质量不错的参考文章了，比如王多鱼的实践DSSM召回（如果对dssm模型原理不熟，建议先阅读这篇文章，再看本文实践部分，本文主要讲实现），总结的非常不错，王多鱼这篇文章DSSM实践是基于浅梦大佬开源deepmatch包实现的，但是在推荐系统实践中如果直接调用别人的模型包会遇到诸多不便，需要在自己业务场景中做finetune；实际生产中，模型所用到的特征往往都是稀疏的，多值变长的，对有些特征我们还想让它们共享embedding，说到这里，我要非常感谢石塔西的这篇文章用TensorFlow实现支持多值、稀疏、共享权重的DeepFM，从这篇文章中，我得到很多启发;本文下面介绍的主要是自己从各位大佬那学习到的知识总结，并无什么创新点，希望对一些刚入坑的童鞋们有所帮助。好了，那我们开始，Talk is cheap, Show me the code.

　　虽然本文不讲模型原理，但是有两Tricks，还是值得提下，亲测有效，这两tricks在下文实现均有体现。

**1　对user及item embedding向量 L2标准化**

$$u(x,\theta) \leftarrow u(x,\theta)/||u(x,\theta)||_2 , \quad v(x,\theta) \leftarrow v(x,\theta)/||v(x,\theta)||_2$$

Emebdding标准化可以加速模型训练和提升检索效果。

**2　增强Softmax效果**

通过引入超参数 $\tau$ 来增强softmax每个逻辑值的输出：

$$s(x,y) = < u(x,\theta), v(y,\theta) > /\tau$$

微调超参数 $\tau$ 可以最大化召回率或精确率。

<div align="center">

| 1 |
|---|

</div>

<div align="center">

| 数据预处理 |
|---|

</div>

为了逼近真实推荐系统场景的数据处理，这里人为构造部分实际生产数据样例作为演示；

| | act | client_id | post_id | client_type | follow_topic_id | all_topic_fav_7 | topic_id | read_post_id |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 28401 | 39647119 | 0 | 572,92,62,37,35,34,33,32,31,30,29,68,67,65,24,... | 502:0.3443,278:0.0868,177:0.0719,1:0.497 | 135 | 39588887,39599018,39576294,39553374,39630091 |
| 1 | 1 | 28401 | 39645671 | 0 | 572,92,62,37,35,34,33,32,31,30,29,68,67,65,24,... | 502:0.3443,278:0.0868,177:0.0719,1:0.497 | 1 | 39588887,39599018,39576294,39553374,39630091 |
| 2 | 0 | 28401 | 39643183 | 0 | 572,92,62,37,35,34,33,32,31,30,29,68,67,65,24,... | 502:0.3443,278:0.0868,177:0.0719,1:0.497 | 3 | 39588887,39599018,39576294,39553374,39630091 |
| 3 | 0 | 28401 | 39629847 | 0 | 572,92,62,37,35,34,33,32,31,30,29,68,67,65,24,... | 502:0.3443,278:0.0868,177:0.0719,1:0.497 | 4 | 39588887,39599018,39576294,39553374,39630091 |
| 4 | 1 | 28401 | 39613538 | 0 | 572,92,62,37,35,34,33,32,31,30,29,68,67,65,24,... | 502:0.3443,278:0.0868,177:0.0719,1:0.497 | 278 | 39588887,39599018,39576294,39553374,39630091 |

样本

字段介绍：

**act**：为label数据 1:正样本，0：负样本

**client_id**: 用户id

**post_id**：物料item id 这里称为post_id

**client_type**:用户客户端类型

**follow_topic_id**: 用户关注话题分类id

**all_topic_fav_7**: 用户画像特征，用户最近7天对话题偏爱度刻画，kv键值对形式

**topic_id**: 物料所属的话题

**read_post_id**:用户最近阅读的物料id

## 预训练item embedding 向量

```
ITEM_EMBEDDING
```

```
<tf.Tensor: shape=(112396, 768), dtype=float32, numpy=
array([[ 0.       , 0.       , 0.       , ..., 0.       , 0.       ,
         0.       ],
       [ 0.999357, 0.999952, 0.997325, ..., -0.997149, -0.999041,
        -0.149874],
       [ 0.999701, 0.999996, 0.475305, ..., -0.996938, -0.999869,
         0.479093],
       ...,
       [ 0.999113, 0.999849, 0.961474, ..., -0.997135, -0.997072,
         0.113319],
       [ 0.999691, 0.999915, 0.978405, ..., -0.999732, -0.996863,
         0.813737],
       [ 0.999831, 0.999961, 0.999492, ..., -0.99991 , -0.999621,
        -0.260133]], dtype=float32)>
```

这里会有为每个item预训练生成一个embedding向量，存到embedding矩阵中，idx=0行，为一个默认值，当一个item因某些原因未生成其embedding向量，则用默认值0替代。

## 定义参数类型

　　我们将参数归三种类型单值离散型SparseFeat，如topic_id字段；稠密数值类型DenseFeat，如用户访问时间及用户embedding向量等；多值变长离散特征VarLenSparseFeat，如follow_topic_id或者带权重形式all_topic_fav_7；这里延用deepMatch开源包里定义输入变量方式，需要注意的是，SparseFeat与VarLenSparseFeat类型的特征，如果想共享embedding权重向量，需要指定其与哪个category离散变量特征embedding参数共享，如这里我们想follow_topic_id与all_topic_fav_7里的id embedding与item的topic_id embedding权重共享一套，设置share_embed='topic_id'即可。

```python
1  from collections import namedtuple, OrderedDict
2  import tensorflow as tf
3
4  SparseFeat = namedtuple('SparseFeat', ['name', 'voc_size', 'share_embed','emb
5  DenseFeat = namedtuple('DenseFeat', ['name', 'pre_embed','reduce_type','dim',
6  VarLenSparseFeat = namedtuple('VarLenSparseFeat', ['name', 'voc_size', 'share
7  import tensorflow as tf
8
9  SparseFeat = namedtuple('SparseFeat', ['name', 'voc_size', 'share_embed','emb
10 DenseFeat = namedtuple('DenseFeat', ['name', 'pre_embed','reduce_type','dim',
11 VarLenSparseFeat = namedtuple('VarLenSparseFeat', ['name', 'voc_size', 'share
```

## 定义DSSM输入变量参数

　　除了常见的特征，这里使用用户最近浏览的物料embedding向量的平均作为用户的一个特征即client_embed；我们将follow_topic_id，all_topic_fav7用到的topic_id embedding向量与item的topic_id对应的embedding向量共享，在实际应用中，相近语义的embedding权重共享是很有必要的，大大减少网络训练参数，防止过拟合。

```python
1  feature_columns = [SparseFeat(name="topic_id", voc_size=700, share_embed=None
2                     SparseFeat(name='client_type', voc_size=2, share_embed=No
3                     VarLenSparseFeat(name="follow_topic_id", voc_size=700, sh
4                     VarLenSparseFeat(name="all_topic_fav_7", voc_size=700, sh
5                     DenseFeat(name='item_embed',pre_embed='post_id', reduce_t
6                     DenseFeat(name='client_embed',pre_embed='read_post_id', r
7                     ]
8
```

```
 9
10    # 用户特征及贴子特征
11    user_feature_columns_name = ["follow_topic_id", 'all_topic_fav_7','client_typ
12    item_feature_columns_name = ["topic_id", 'post_type','item_embed',]
13    user_feature_columns = [col for col in feature_columns if col.name in user_fe
14    item_feature_columns = [col for col in feature_columns if col.name in item_fe
```

## 构造训练tf.dataset数据

首先加载预训练 item embedding向量及离散特征vocabulary

```
 1    def get_item_embed(file_names):
 2        item_bert_embed = []
 3        item_id = []
 4        for file in file_names:
 5            with open(file, 'r') as f:
 6                for line in f:
 7                    feature_json = json.loads(line)
 8                    item_bert_embed.append(feature_json['post_id'])
 9                    item_id.append(feature_json['values'])
10
11        item_id2idx = tf.lookup.StaticHashTable(
12            tf.lookup.KeyValueTensorInitializer(
13                keys=item_id,
14                values=range(1, len(item_id)+1),
15                key_dtype=tf.string,
16                value_dtype=tf.int32),
17            default_value=0)
18        item_bert_embed = [[0.0]*768] + item_bert_embed
19        item_embedding = tf.constant(item_bert_embed, dtype=tf.float32)
20        return item_id2idx, item_embedding
21    # 获取item embedding及其查找关系
22    ITEM_ID2IDX, ITEM_EMBEDDING = get_item_embed(file_names)
23
24    # 定义离散特征集合，离散特征vocabulary
25    DICT_CATEGORICAL = {"topic_id": [str(i) for i in range(0, 700)],
26                "client_type": [0,1]
27                }
```

然后，tf.dataset构造

```
DEFAULT_VALUES = [[0],[''],[''],[0.0], [''], [''], [''],['']]
COL_NAME = ['act', 'client_id', 'post_id', 'client_type', 'follow_topic_id',

def _parse_function(example_proto):

    item_feats = tf.io.decode_csv(example_proto, record_defaults=DEFAULT_VAl
    parsed = dict(zip(COL_NAME, item_feats))

    feature_dict = {}
    for feat_col in feature_columns:
        if isinstance(feat_col, VarLenSparseFeat):
            if feat_col.weight_name is not None:
                kvpairs = tf.strings.split([parsed[feat_col.name]], ',').va
                kvpairs = tf.strings.split(kvpairs, ':')
                kvpairs = kvpairs.to_tensor()
                feat_ids, feat_vals = tf.split(kvpairs, num_or_size_splits=2
                feat_vals= tf.strings.to_number(feat_vals, out_type=tf.floa
                feature_dict[feat_col.name] = feat_ids
                feature_dict[feat_col.weight_name] = feat_vals
            else:
                feat_ids = tf.strings.split([parsed[feat_col.name]], ',').va
                feat_ids = tf.reshape(feat_ids, shape=[-1])
                feature_dict[feat_col.name] = feat_ids

        elif isinstance(feat_col, SparseFeat):
            feature_dict[feat_col.name] = parsed[feat_col.name]

        elif isinstance(feat_col, DenseFeat):
            if feat_col.pre_embed is None:
                feature_dict[feat_col.name] = parsed[feat_col.name]
            elif feat_col.reduce_type is not None:
                keys = tf.strings.split(parsed[feat_col.pre_embed], ',')
                emb = tf.nn.embedding_lookup(params=ITEM_EMBEDDING, ids=ITEN
                emb = tf.reduce_mean(emb,axis=0) if feat_col.reduce_type ==
                feature_dict[feat_col.name] = emb
            else:
                emb = tf.nn.embedding_lookup(params=ITEM_EMBEDDING, ids=ITEN
                feature_dict[feat_col.name] = emb
```

```python
39              else:
40                  raise "unknown feature_columns...."
41
42
43      label = parsed['act']
44
45
46      return feature_dict, label
47
48
49  pad_shapes = {}
50  pad_values = {}
51
52  for feat_col in feature_columns:
53      if isinstance(feat_col, VarLenSparseFeat):
54          max_tokens = feat_col.maxlen
55          pad_shapes[feat_col.name] = tf.TensorShape([max_tokens])
56          pad_values[feat_col.name] = ''
57          if feat_col.weight_name is not None:
58              pad_shapes[feat_col.weight_name] = tf.TensorShape([max_tokens])
59              pad_values[feat_col.weight_name] = tf.constant(-1, dtype=tf.floa
60
61  # no need to pad labels
62      elif isinstance(feat_col, SparseFeat):
63          if feat_col.dtype == 'string':
64              pad_shapes[feat_col.name] = tf.TensorShape([])
65              pad_values[feat_col.name] = '9999'
66          else:
67              pad_shapes[feat_col.name] = tf.TensorShape([])
68              pad_values[feat_col.name] = 0.0
69      elif isinstance(feat_col, DenseFeat):
70          if feat_col.pre_embed is None:
71              pad_shapes[feat_col.name] = tf.TensorShape([])
72              pad_values[feat_col.name] = 0.0
73          else:
74              pad_shapes[feat_col.name] = tf.TensorShape([feat_col.dim])
75              pad_values[feat_col.name] = 0.0
76
77
78  pad_shapes = (pad_shapes, (tf.TensorShape([])))
```

```python
79    pad_values = (pad_values, (tf.constant(0, dtype=tf.int32)))
80
81
82    filenames= tf.data.Dataset.list_files([
83    '/recall_user_item_act.csv'
84    ])
85    dataset = filenames.flat_map(
86            lambda filepath: tf.data.TextLineDataset(filepath).skip(1))
87
88    batch_size = 1024
89    dataset = dataset.map(_parse_function, num_parallel_calls=60)
90    dataset = dataset.repeat()
91    dataset = dataset.shuffle(buffer_size = batch_size*2) # 在缓冲区中随机打乱数据
92    dataset = dataset.padded_batch(batch_size = batch_size,
93                                      padded_shapes = pad_shapes,
94                                      padding_values = pad_values) # 每1024条数据为一
95    dataset = dataset.prefetch(buffer_size=tf.data.experimental.AUTOTUNE)
96
97    # 验证集
98    filenames_val= tf.data.Dataset.list_files(['/recall_user_item_act_val.csv'])
99    dataset_val = filenames_val.flat_map(
100            lambda filepath: tf.data.TextLineDataset(filepath).skip(1))
101
102   val_batch_size = 1024
103   dataset_val = dataset_val.map(_parse_function, num_parallel_calls=60)
104   dataset_val = dataset_val.padded_batch(batch_size = val_batch_size,
105                                      padded_shapes = pad_shapes,
106                                      padding_values = pad_values) # 每1024条数据为一
107   dataset_val = dataset_val.prefetch(buffer_size=tf.data.experimental.AUTOTUNE
```

经过上述逻辑代码预处理后，原始样本csv文件中数据格式已经转化为如下的形式（这里拿batch_size=1 举例），kv形式的特征被拆分为两个Input输入变量一个是category离散ID（如all_topic_fav_7），一个是其对应的weight（如all_topic_fav_7_weight），他们最终被输入到tf.nn.embedding_lookup_sparse(self.embedding,sp_ids=idx,　　　　　　sp_weights=val, combiner='sum') 这个api 对应的sp_ids，sp_weights参数中去。

```python
1    # next(iter(dataset))
2    ({'topic_id': <tf.Tensor: shape=(1,), dtype=string, numpy=array([b'278'], dty
```

```
3     'client_type': <tf.Tensor: shape=(1,), dtype=float32, numpy=array([0.], dty
4     'follow_topic_id': <tf.Tensor: shape=(1, 20), dtype=string, numpy=
5     array([[b'572', b'92', b'62', b'37', b'35', b'34', b'33', b'32', b'31',
6            b'30', b'29', b'68', b'67', b'65', b'24', b'20', b'16', b'15',
7            b'13', b'12']], dtype=object)>,
8     'all_topic_fav_7': <tf.Tensor: shape=(1, 5), dtype=string, numpy=array([[b'
9     'all_topic_fav_7_weight': <tf.Tensor: shape=(1, 5), dtype=float32, numpy=ar
10    'item_embed': <tf.Tensor: shape=(1, 768), dtype=float32, numpy=
11    array([[ 0.999586,  0.999861,  0.995566,  0.892292,  0.848516,  0.815888,
12           -0.860286, -0.871219,  0.982316, -0.999692,  0.999998,  0.999589,
13                                     ......
14           -0.943752,  0.999957, -0.990231,  0.999377, -0.997795,  0.999498,
15           -0.995729,  0.701236,  0.991473,  0.946505, -0.996337,  0.999991,
16            0.991516, -0.997269, -0.993377, -0.9964  , -0.99972 ,  0.880781]],
17         dtype=float32)>,
18    'client_embed': <tf.Tensor: shape=(1, 768), dtype=float32, numpy=
19    array([[ 0.79698  ,  0.7999152 ,  0.78845704,  0.6598178 ,  0.59617054,
20            0.5318628 , -0.5754676 , -0.7469004 ,  0.78916025, -0.7958456 ,
21                                     ......
22            0.7989754 , -0.7971929 , -0.0165708 ,  0.7924882 ,  0.73336124,
23           -0.794997  ,  0.7999618 ,  0.7634414 , -0.792517  , -0.762231   ,
24           -0.7960204 , -0.7998554 ,  0.37363502]], dtype=float32)>},
25    <tf.Tensor: shape=(1,), dtype=int32, numpy=array([1], dtype=int32)>)
```

2

自定义模型层

```
1     # 离散多值查找表 转稀疏SparseTensor >> EncodeMultiEmbedding >>tf.nn.embedding_
2     class SparseVocabLayer(Layer):
3         def __init__(self, keys, **kwargs):
4             super(SparseVocabLayer, self).__init__(**kwargs)
5             vals = tf.range(1, len(keys) + 1)
6             vals = tf.constant(vals, dtype=tf.int32)
```

```python
 7          keys = tf.constant(keys)
 8          self.table = tf.lookup.StaticHashTable(
 9              tf.lookup.KeyValueTensorInitializer(keys, vals), 0)
10
11      def call(self, inputs):
12          input_idx = tf.where(tf.not_equal(inputs, ''))
13          input_sparse = tf.SparseTensor(input_idx, tf.gather_nd(inputs, input
14          return tf.SparseTensor(indices=input_sparse.indices,
15                                 values=self.table.lookup(input_sparse.values)
16                                 dense_shape=input_sparse.dense_shape)
17
18  # 自定义Embedding层，初始化时，需要传入预先定义好的embedding矩阵，好处可以共享embe
19  class EncodeMultiEmbedding(Layer):
20      def __init__(self, embedding, has_weight=False, **kwargs):
21
22          super(EncodeMultiEmbedding, self).__init__(**kwargs)
23          self.has_weight = has_weight
24          self.embedding = embedding
25
26
27      def build(self, input_shape):
28          super(EncodeMultiEmbedding, self).build(input_shape)
29
30      def call(self, inputs):
31          if self.has_weight:
32              idx, val = inputs
33              combiner_embed = tf.nn.embedding_lookup_sparse(self.embedding,sp
34          else:
35              idx = inputs
36              combiner_embed = tf.nn.embedding_lookup_sparse(self.embedding,sp
37          return tf.expand_dims(combiner_embed, 1)
38
39      def get_config(self):
40          config = super(EncodeMultiEmbedding, self).get_config()
41          config.update({'has_weight': self.has_weight})
42          return config
43
44  # 稠密权重转稀疏格式输入到tf.nn.embedding_lookup_sparse的sp_weights参数中
45  class Dense2SparseTensor(Layer):
46      def __init__(self):
```

```python
        super(Dense2SparseTensor, self).__init__()

    def call(self, dense_tensor):
        weight_idx = tf.where(tf.not_equal(dense_tensor, tf.constant(-1, dty
        weight_sparse = tf.SparseTensor(weight_idx, tf.gather_nd(dense_tens
        return weight_sparse

    def get_config(self):
        config = super(Dense2SparseTensor, self).get_config()
        return config


# 自定义dnese层含BN， dropout
class CustomDense(Layer):
    def __init__(self, units=32, activation='tanh', dropout_rate =0, use_bn=
        self.units = units
        self.activation = activation
        self.dropout_rate = dropout_rate
        self.use_bn = use_bn
        self.seed = seed
        self.tag_name = tag_name

        super(CustomDense, self).__init__(**kwargs)

    #build方法一般定义Layer需要被训练的参数。
    def build(self, input_shape):
        self.weight = self.add_weight(shape=(input_shape[-1], self.units),
                                initializer='random_normal',
                                trainable=True,
                                name='kernel_' + self.tag_name)
        self.bias = self.add_weight(shape=(self.units,),
                                initializer='random_normal',
                                trainable=True,
                                name='bias_' + self.tag_name)

        if self.use_bn:
            self.bn_layers = tf.keras.layers.BatchNormalization()

        self.dropout_layers = tf.keras.layers.Dropout(self.dropout_rate)
        self.activation_layers = tf.keras.layers.Activation(self.activation
```

```python
 87
 88            super(CustomDense,self).build(input_shape) # 相当于设置self.built = 1

 90        #call方法一般定义正向传播运算逻辑，__call__方法调用了它。
 91        def call(self, inputs, training = None, **kwargs):
 92            fc = tf.matmul(inputs, self.weight) + self.bias
 93            if self.use_bn:
 94                fc = self.bn_layers(fc)
 95            out_fc = self.activation_layers(fc)

 97            return out_fc

 99        #如果要让自定义的Layer通过Functional API 组合成模型时可以序列化，需要自定义get
100        def get_config(self):
101            config = super(CustomDense, self).get_config()
102            config.update({'units': self.units, 'activation': self.activation,
103                           'dropout_rate': self.dropout_rate, 'seed': self.seed,
104            return config



# cos 相似度计算层
class Similarity(Layer):

    def __init__(self, gamma=1, axis=-1, type_sim='cos', **kwargs):
        self.gamma = gamma
        self.axis = axis
        self.type_sim = type_sim
        super(Similarity, self).__init__(**kwargs)

    def build(self, input_shape):
        # Be sure to call this somewhere!
        super(Similarity, self).build(input_shape)

    def call(self, inputs, **kwargs):
        query, candidate = inputs
        if self.type_sim == "cos":
            query_norm = tf.norm(query, axis=self.axis)
            candidate_norm = tf.norm(candidate, axis=self.axis)
        cosine_score = tf.reduce_sum(tf.multiply(query, candidate), -1)
        cosine_score = tf.divide(cosine_score, query_norm * candidate_norm
```

```python
127             cosine_score = tf.clip_by_value(cosine_score, -1, 1.0) * self.gamma
128             return tf.expand_dims(cosine_score, 1)
129
130         def compute_output_shape(self, input_shape):
131             return (None, 1)
132
133         def get_config(self, ):
134             config = {'gamma': self.gamma, 'axis': self.axis, 'type': self.type_
135             base_config = super(Similarity, self).get_config()
136             return base_config.uptate(config)
137
138
139     # 自定损失函数，加权交叉熵损失
140     class WeightedBinaryCrossEntropy(tf.keras.losses.Loss):
141         """
142         Args:
143           pos_weight: Scalar to affect the positive labels of the loss function
144           weight: Scalar to affect the entirety of the loss function.
145           from_logits: Whether to compute loss from logits or the probability.
146           reduction: Type of tf.keras.losses.Reduction to apply to loss.
147           name: Name of the loss function.
148         """
149
150         def __init__(self, pos_weight=1.2, from_logits=False,
151                      reduction=tf.keras.losses.Reduction.AUTO,
152                      name='weighted_binary_crossentropy'):
153             super().__init__(reduction=reduction, name=name)
154             self.pos_weight = pos_weight
155             self.from_logits = from_logits
156
157         def call(self, y_true, y_pred):
158             y_true = tf.cast(y_true, tf.float32)
159             ce = tf.losses.binary_crossentropy(
160                 y_true, y_pred, from_logits=self.from_logits)[:, None]
161             ce = ce * (1 - y_true) + self.pos_weight * ce * (y_true)
162     #         ce =tf.nn.weighted_cross_entropy_with_logits(
163     #             y_true, y_pred, self.pos_weight, name=None
164     #         )
165
166             return ce
```

```
167
168    def get_config(self, ):
169        config = {'pos_weight': self.pos_weight, 'from_logits': self.from_l
170        base_config = super(WeightedBinaryCrossEntropy, self).get_config()
171        return base_config.uptate(config)
```

3

定义输入及共享层帮助函数

```
1   # 定义model输入特征
2   def build_input_features(features_columns, prefix=''):
3       input_features = OrderedDict()
4       for feat_col in features_columns:
5           if isinstance(feat_col, DenseFeat):
6               if feat_col.pre_embed is None:
7                   input_features[feat_col.name] = Input([1], name=feat_col.nam
8               else:
9                   input_features[feat_col.name] = Input([feat_col.dim], name=1
10          elif isinstance(feat_col, SparseFeat):
11              if feat_col.dtype == 'string':
12                  input_features[feat_col.name] = Input([None], name=feat_col
13              else:
14                  input_features[feat_col.name] = Input([1], name=feat_col.nam
15          elif isinstance(feat_col, VarLenSparseFeat):
16              input_features[feat_col.name] = Input([None], name=feat_col.name
17              if feat_col.weight_name is not None:
18                  input_features[feat_col.weight_name] = Input([None], name=fe
19          else:
20              raise TypeError("Invalid feature column in build_input_features
21
22      return input_features
23
24  # 构造自定义embedding层matrix
```

```python
def build_embedding_matrix(features_columns):
    embedding_matrix = {}
    for feat_col in features_columns:
        if isinstance(feat_col, SparseFeat) or isinstance(feat_col, VarLenS
            if feat_col.dtype == 'string':
                vocab_name = feat_col.share_embed if feat_col.share_embed e
                vocab_size = feat_col.voc_size
                embed_dim = feat_col.embed_dim
                if vocab_name not in embedding_matrix:
                    embedding_matrix[vocab_name] = tf.Variable(initial_value

    return embedding_matrix

# 构造自定义 embedding层
def build_embedding_dict(features_columns, embedding_matrix):
    embedding_dict = {}
    for feat_col in features_columns:
        if isinstance(feat_col, SparseFeat):
            if feat_col.dtype == 'string':
                vocab_name = feat_col.share_embed if feat_col.share_embed e
                embedding_dict[feat_col.name] = EncodeMultiEmbedding(embedd
        elif isinstance(feat_col, VarLenSparseFeat):
            vocab_name = feat_col.share_embed if feat_col.share_embed else
            if feat_col.weight_name is not None:
                embedding_dict[feat_col.name] = EncodeMultiEmbedding(embedd
            else:
                embedding_dict[feat_col.name] = EncodeMultiEmbedding(embedd

    return embedding_dict


# dense 与 embedding特征输入
def input_from_feature_columns(features, features_columns, embedding_dict):
    sparse_embedding_list = []
    dense_value_list = []

    for feat_col in features_columns:
        if isinstance(feat_col, SparseFeat) or isinstance(feat_col, VarLenS
            if feat_col.dtype == 'string':
                vocab_name = feat_col.share_embed if feat_col.share_embed e
```

```
65                    keys = DICT_CATEGORICAL[vocab_name]
66                    _input_sparse = SparseVocabLayer(keys)(features[feat_col.nar
67
68            if isinstance(feat_col, SparseFeat):
69                if feat_col.dtype == 'string':
70                    _embed = embedding_dict[feat_col.name](_input_sparse)
71                else:
72                    _embed = Embedding(feat_col.voc_size+1, feat_col.embed_dim,
73                                       embeddings_regularizer=tf.keras.regularizers
74                sparse_embedding_list.append(_embed)
75            elif isinstance(feat_col, VarLenSparseFeat):
76                if feat_col.weight_name is not None:
77                    _weight_sparse = Dense2SparseTensor()(features[feat_col.weig
78                    _embed = embedding_dict[feat_col.name]([_input_sparse, _weig
79
80                else:
81                    _embed = embedding_dict[feat_col.name](_input_sparse)
82                sparse_embedding_list.append(_embed)
83
84            elif isinstance(feat_col, DenseFeat):
85                dense_value_list.append(features[feat_col.name])
86
87            else:
88                raise TypeError("Invalid feature column in input_from_feature_co
89
90        return sparse_embedding_list, dense_value_list
91
92
93  def concat_func(inputs, axis=-1):
94      if len(inputs) == 1:
95          return inputs[0]
96      else:
97          return Concatenate(axis=axis)(inputs)
98
99  def combined_dnn_input(sparse_embedding_list, dense_value_list):
100     if len(sparse_embedding_list) > 0 and len(dense_value_list) > 0:
101         sparse_dnn_input = Flatten()(concat_func(sparse_embedding_list))
102         dense_dnn_input = Flatten()(concat_func(dense_value_list))
103         return concat_func([sparse_dnn_input, dense_dnn_input])
104     elif len(sparse_embedding_list) > 0:
```

```python
105        return Flatten()(concat_func(sparse_embedding_list))
106    elif len(dense_value_list) > 0:
107        return Flatten()(concat_func(dense_value_list))
108    else:
109        raise "dnn_feature_columns can not be empty list"
```

4

搭建DSSM模型

```python
1  def DSSM(
2      user_feature_columns,
3      item_feature_columns,
4      user_dnn_hidden_units=(256, 256, 128),
5      item_dnn_hidden_units=(256, 256,  128),
6      user_dnn_dropout=(0, 0, 0),
7      item_dnn_dropout=(0, 0, 0),
8      out_dnn_activation='tanh',
9      gamma=1.2,
10     dnn_use_bn=False,
11     seed=1024,
12     metric='cos'):
13
14     """
15     Instantiates the Deep Structured Semantic Model architecture.
16     Args:
17         user_feature_columns: A list containing user's features used by the m
18         item_feature_columns: A list containing item's features used by the m
19         user_dnn_hidden_units: tuple,tuple of positive integer , the layer nu
20         item_dnn_hidden_units: tuple,tuple of positive integer, the layer num
21         out_dnn_activation: Activation function to use in deep net
22         dnn_use_bn: bool. Whether use BatchNormalization before activation or
23         user_dnn_dropout: tuple of float in [0,1), the probability we will dr
24         item_dnn_dropout: tuple of float in [0,1), the probability we will dr
```

```python
25          seed: integer ,to use as random seed.
26          gamma: A useful hyperparameter for Similarity layer
27          metric: str, "cos" for  cosine
28      return: A TF Keras model instance.
29      """
30      features_columns = user_feature_columns + item_feature_columns
31      # 构建 embedding_dict
32      embedding_matrix = build_embedding_matrix(features_columns)
33      embedding_dict = build_embedding_dict(features_columns, embedding_matrix)
34
35      # user 特征 处理
36      user_features = build_input_features(user_feature_columns)
37      user_inputs_list = list(user_features.values())
38      user_sparse_embedding_list, user_dense_value_list = input_from_feature_cc
39
40      user_dnn_input = combined_dnn_input(user_sparse_embedding_list, user_dens
41
42      # item 特征 处理
43      item_features = build_input_features(item_feature_columns)
44      item_inputs_list = list(item_features.values())
45      item_sparse_embedding_list, item_dense_value_list = input_from_feature_cc
46
47      item_dnn_input = combined_dnn_input(item_sparse_embedding_list, item_dens
48
49
50      # user tower
51      for i in range(len(user_dnn_hidden_units)):
52          if i == len(user_dnn_hidden_units) - 1:
53              user_dnn_out = CustomDense(units=user_dnn_hidden_units[i],dropout
54                                        use_bn=dnn_use_bn,activation=out_dnn_a
55              break
56          user_dnn_input = CustomDense(units=user_dnn_hidden_units[i],dropout_r
57                                      use_bn=dnn_use_bn,activation='relu', nam
58
59
60      # item tower
61      for i in range(len(item_dnn_hidden_units)):
62          if i == len(item_dnn_hidden_units) - 1:
63              item_dnn_out = CustomDense(units=item_dnn_hidden_units[i],dropout
64                                        use_bn=dnn_use_bn, activation=out_dnn_
```

```
65              break
66          item_dnn_input = CustomDense(units=item_dnn_hidden_units[i],dropout_r
67                              use_bn=dnn_use_bn,activation='relu', nam
68
69
70      score = Similarity(type_sim=metric,gamma=gamma)([user_dnn_out, item_dnn_o
71      output = tf.keras.layers.Activation("sigmoid", name="dssm_out")(score)
72  #    score = Multiply()([user_dnn_out, item_dnn_out])
73  #    output = Dense(1, activation="sigmoid",name="dssm_out")(score)
74
75      model = Model(inputs=user_inputs_list + item_inputs_list, outputs=output)
76      model.__setattr__("user_input", user_inputs_list)
77      model.__setattr__("item_input", item_inputs_list)
78      model.__setattr__("user_embedding", user_dnn_out)
79      model.__setattr__("item_embedding", item_dnn_out)
80
81      return model
```

<div align="center">5</div>

<div align="center">训练及保存模型</div>

## 训练模型

```
1   model= DSSM(
2       user_feature_columns,
3       item_feature_columns,
4       user_dnn_hidden_units=(256, 256, 128),
5       item_dnn_hidden_units=(256, 256, 128),
6       user_dnn_dropout=(0, 0, 0),
7       item_dnn_dropout=(0, 0, 0),
8       out_dnn_activation='tanh',
9       gamma=1,
10      dnn_use_bn=False,
11      seed=1024,
```

```
12          metric='cos')
13
14  model.compile(optimizer='adagrad',
15                loss={"dssm_out": WeightedBinaryCrossEntropy(),
16                      },
17                loss_weights=[1.0,],
18                metrics={"dssm_out": [tf.keras.metrics.AUC(name='auc')]}
19                )
20
21
22  log_dir = '/mywork/tensorboardshare/logs/' + datetime.datetime.now().strftime
23  tbCallBack = TensorBoard(log_dir=log_dir,  # log 目录
24                histogram_freq=0,  # 按照何等频率（epoch）来计算直方图，0为不计算
25                write_graph=True,  # 是否存储网络结构图
26                write_images=True,# 是否可视化参数
27                update_freq='epoch',
28                embeddings_freq=0,
29                embeddings_layer_names=None,
30                embeddings_metadata=None,
31                    profile_batch = 40)
32
33  #
34  #
35  total_train_sample =  115930
36  total_test_sample =   1181
37  train_steps_per_epoch=np.floor(total_train_sample/batch_size).astype(np.int32
38  test_steps_per_epoch = np.ceil(total_test_sample/val_batch_size).astype(np.in
39  history_loss = model.fit(dataset, epochs=1,
40          steps_per_epoch=train_steps_per_epoch,
41          validation_data=dataset_val, validation_steps=test_steps_per_epoch,
42          verbose=1,callbacks=[tbCallBack])
```
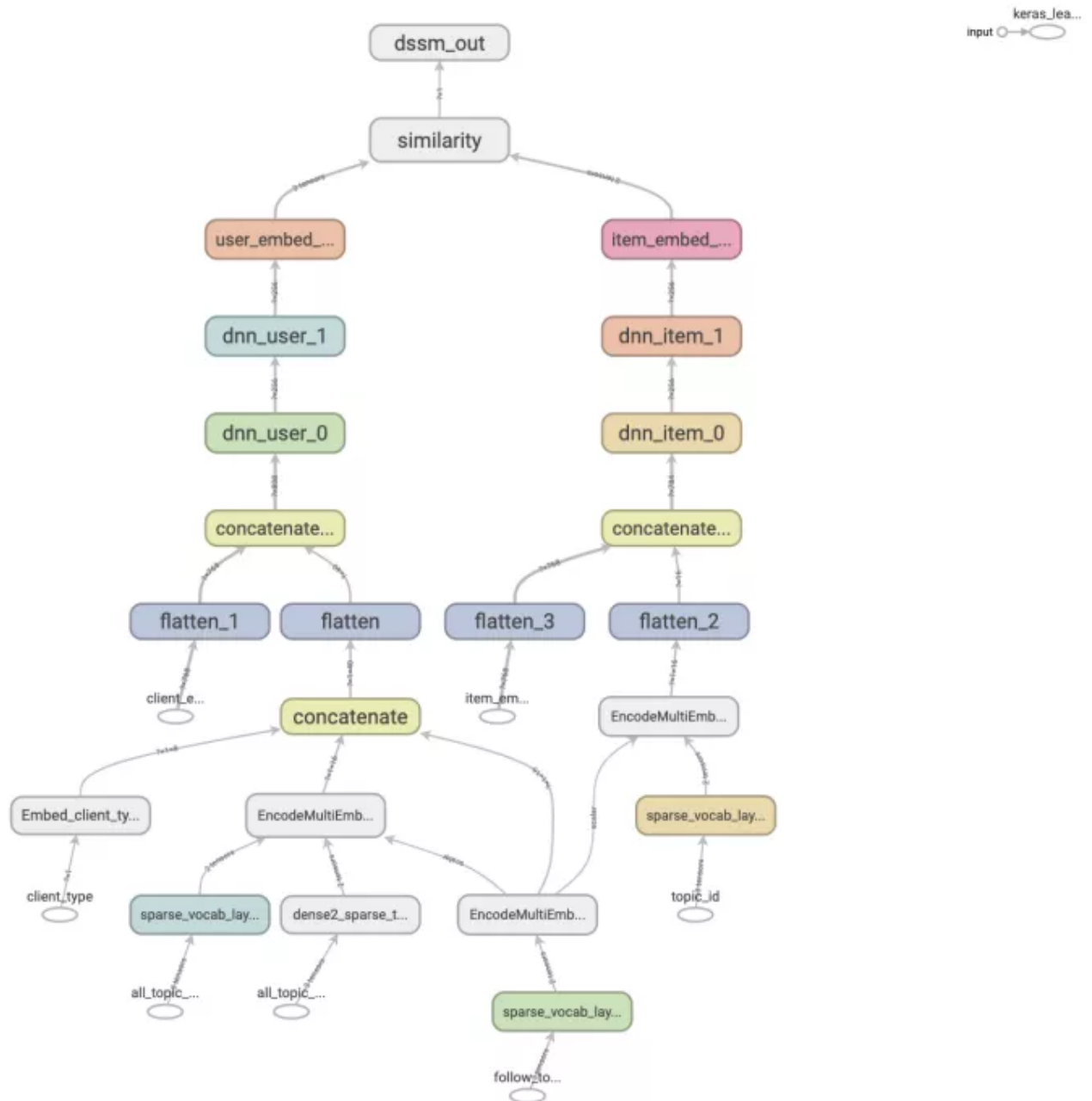
## 模型结构summary

## 保存模型

```
1   # 用户塔 item塔定义
2   user_embedding_model = Model(inputs=model.user_input, outputs=model.user_embed
3   item_embedding_model = Model(inputs=model.item_input, outputs=model.item_embed
4   # 保存
5   tf.keras.models.save_model(user_embedding_model,"/Recall/DSSM/models/dssmUser/
6   tf.keras.models.save_model(item_embedding_model,"/Recall/DSSM/models/dssmItem/
```

## 获取user embedding 及item embedding

```python
user_query = {'all_topic_fav_7': np.array([['294', '88', '60', '1']]),
  'all_topic_fav_7_weight':np.array([[ 0.0897,  0.2464,  0.0928,  0.5711,]]),
  'follow_topic_id': np.array([['75', '73', '74', '92', '62', '37', '35', '34
  'client_type': np.array([0.]),
    'client_embed': np.array([[-9.936600e-02,  2.752400e-01, -4.314620e-01,
          -5.263000e-02, -4.490300e-01, -3.641180e-01, -3.545410e-01,
          -2.315470e-01,  4.641480e-01,  3.965120e-01, -1.670170e-01,
          -5.480000e-03, -1.646790e-01,  2.522832e+00, -2.946590e-01,
                                ......
          -1.151946e+00, -4.008270e-01,  1.521650e-01, -3.524520e-01,
           4.836160e-01, -1.190920e-01,  5.792700e-02, -6.148070e-01,
          -7.182930e-01, -1.351920e-01,  2.048980e-01, -1.259220e-01]])}

item_query = {
  'topic_id': np.array(['1']),
  'item_embed': np.array([[-9.936600e-02,  2.752400e-01, -4.314620e-01,  3.39
          -5.263000e-02, -4.490300e-01, -3.641180e-01, -3.545410e-01,
          -2.315470e-01,  4.641480e-01,  3.965120e-01, -1.670170e-01,
                                ......
          -1.151946e+00, -4.008270e-01,  1.521650e-01, -3.524520e-01,
           4.836160e-01, -1.190920e-01,  5.792700e-02, -6.148070e-01,
          -7.182930e-01, -1.351920e-01,  2.048980e-01, -1.259220e-01]]),
}

user_embs = user_embedding_model.predict(user_query)
item_embs = item_embedding_model.predict(item_query)

# 结果:
# user_embs:
# array([[ 0.80766946,  0.13907856, -0.37779272,  0.53268254, -0.3095821 ,
#           0.2213103 , -0.24618168, -0.7127088 ,  0.4502724 ,  0.4282374 ,
#          -0.36033005,  0.43310016, -0.29158285,  0.8743557 ,  0.5113318 ,
#           0.26994514, -0.35604447,  0.33559784, -0.28052363,  0.38596702,
#           0.5038488 , -0.32811972, -0.5471834 , -0.07594685,  0.7006799 ,
#          -0.24201767,  0.31005877, -0.06173763, -0.28473467,  0.61975694,
#                                ......
#          -0.714099  , -0.5384026 ,  0.38787717, -0.4263588 ,  0.30690318,
#           0.24047776, -0.01420124,  0.15475503,  0.77783686, -0.43002903,
#           0.52561694,  0.37806144,  0.18955356, -0.37184635,  0.5181224 ,
#          -0.18585253,  0.05573007, -0.38589332, -0.7673693 , -0.25266737,
```

```
41   #          0.51427466,  0.47647673,  0.47982445]], dtype=float32)
42   # item_embs：
43   # array([[-6.9417924e-01, -3.9942840e-01,  7.2445291e-01, -5.8977932e-01,
44   #          -5.8792406e-01,  5.3883100e-01, -7.8469634e-01,  6.8996024e-01,
45   #          -7.6087400e-02, -4.4855604e-01,  8.4910756e-01, -4.7288817e-01,
46   #          -9.0812451e-01, -4.0452164e-01,  8.8695991e-01, -7.9177713e-01,
47                                    ......
48   #          -9.7515762e-01, -5.2411711e-01,  9.2708725e-01, -1.3903661e-01,
49   #           7.8691095e-01, -8.0726832e-01, -7.3851186e-01,  2.7774110e-01,
50   #          -4.1870885e-02,  4.7335419e-01,  3.4424815e-01, -5.8394599e-01]],
51   #       dtype=float32)
```

6

线上Serving

　　我们这里向量召回检索框架用的是Milvus，用户的UE是线上实时获取的，item的embedding是异步获取存到Milvus平台上。

**参考文献▼**

王多鱼：实践DSSM召回模型

石塔西：用TensorFlow实现支持多值、稀疏、共享权重的DeepFM

https://github.com/shenweichen/