# 16.attention-is-all-you-need

September 29, 2021

```
[1]: import sys
     import warnings

     if not sys.warnoptions:
         warnings.simplefilter('ignore')
```

```
[2]: import tensorflow as tf
     import numpy as np
     import matplotlib.pyplot as plt
     import seaborn as sns
     import pandas as pd
     from sklearn.preprocessing import MinMaxScaler
     from datetime import datetime
     from datetime import timedelta
     from tqdm import tqdm
     sns.set()
     tf.compat.v1.random.set_random_seed(1234)
```

```
[3]: df = pd.read_csv('../dataset/GOOG-year.csv')
     df.head()
```

```
[3]:         Date        Open        High         Low       Close   Adj Close  \
     0  2016-11-02  778.200012  781.650024  763.450012  768.700012  768.700012
     1  2016-11-03  767.250000  769.950012  759.030029  762.130005  762.130005
     2  2016-11-04  750.659973  770.359985  750.560974  762.020020  762.020020
     3  2016-11-07  774.500000  785.190002  772.549988  782.520020  782.520020
     4  2016-11-08  783.400024  795.632996  780.190002  790.510010  790.510010

          Volume
     0   1872400
     1   1943200
     2   2134800
     3   1585100
     4   1350800
```

```
[4]: minmax = MinMaxScaler().fit(df.iloc[:, 4:5].astype('float32')) # Close index
     df_log = minmax.transform(df.iloc[:, 4:5].astype('float32')) # Close index
```

```
df_log = pd.DataFrame(df_log)
df_log.head()
```

[4]:           0
     0  0.112708
     1  0.090008
     2  0.089628
     3  0.160459
     4  0.188066

## 0.1 Split train and test

I will cut the dataset to train and test datasets,

1. Train dataset derived from starting timestamp until last 30 days
2. Test dataset derived from last 30 days until end of the dataset

So we will let the model do forecasting based on last 30 days, and we will going to repeat the experiment for 10 times. You can increase it locally if you want, and tuning parameters will help you by a lot.

```
[5]: test_size = 30
     simulation_size = 10

     df_train = df_log.iloc[:-test_size]
     df_test = df_log.iloc[-test_size:]
     df.shape, df_train.shape, df_test.shape
```

[5]: ((252, 7), (222, 1), (30, 1))

```
[6]: def layer_norm(inputs, epsilon=1e-8):
         mean, variance = tf.nn.moments(inputs, [-1], keep_dims=True)
         normalized = (inputs - mean) / (tf.sqrt(variance + epsilon))

         params_shape = inputs.get_shape()[-1:]
         gamma = tf.get_variable('gamma', params_shape, tf.float32, tf.
     ↪ones_initializer())
         beta = tf.get_variable('beta', params_shape, tf.float32, tf.
     ↪zeros_initializer())

         outputs = gamma * normalized + beta
         return outputs

     def multihead_attn(queries, keys, q_masks, k_masks, future_binding, num_units,␣
     ↪num_heads):

         T_q = tf.shape(queries)[1]
         T_k = tf.shape(keys)[1]
```

```python
    Q = tf.layers.dense(queries, num_units, name='Q')
    K_V = tf.layers.dense(keys, 2*num_units, name='K_V')
    K, V = tf.split(K_V, 2, -1)

    Q_ = tf.concat(tf.split(Q, num_heads, axis=2), axis=0)
    K_ = tf.concat(tf.split(K, num_heads, axis=2), axis=0)
    V_ = tf.concat(tf.split(V, num_heads, axis=2), axis=0)

    align = tf.matmul(Q_, tf.transpose(K_, [0,2,1]))
    align = align / np.sqrt(K_.get_shape().as_list()[-1])

    paddings = tf.fill(tf.shape(align), float('-inf'))

    key_masks = k_masks
    key_masks = tf.tile(key_masks, [num_heads, 1])
    key_masks = tf.tile(tf.expand_dims(key_masks, 1), [1, T_q, 1])
    align = tf.where(tf.equal(key_masks, 0), paddings, align)

    if future_binding:
        lower_tri = tf.ones([T_q, T_k])
        lower_tri = tf.linalg.LinearOperatorLowerTriangular(lower_tri).
→to_dense()
        masks = tf.tile(tf.expand_dims(lower_tri,0), [tf.shape(align)[0], 1, 1])
        align = tf.where(tf.equal(masks, 0), paddings, align)

    align = tf.nn.softmax(align)
    query_masks = tf.to_float(q_masks)
    query_masks = tf.tile(query_masks, [num_heads, 1])
    query_masks = tf.tile(tf.expand_dims(query_masks, -1), [1, 1, T_k])
    align *= query_masks

    outputs = tf.matmul(align, V_)
    outputs = tf.concat(tf.split(outputs, num_heads, axis=0), axis=2)
    outputs += queries
    outputs = layer_norm(outputs)
    return outputs


def pointwise_feedforward(inputs, hidden_units, activation=None):
    outputs = tf.layers.dense(inputs, 4*hidden_units, activation=activation)
    outputs = tf.layers.dense(outputs, hidden_units, activation=None)
    outputs += inputs
    outputs = layer_norm(outputs)
    return outputs
```

```python
def learned_position_encoding(inputs, mask, embed_dim):
    T = tf.shape(inputs)[1]
    outputs = tf.range(tf.shape(inputs)[1])                # (T_q)
    outputs = tf.expand_dims(outputs, 0)                   # (1, T_q)
    outputs = tf.tile(outputs, [tf.shape(inputs)[0], 1])   # (N, T_q)
    outputs = embed_seq(outputs, T, embed_dim, zero_pad=False, scale=False)
    return tf.expand_dims(tf.to_float(mask), -1) * outputs


def sinusoidal_position_encoding(inputs, mask, repr_dim):
    T = tf.shape(inputs)[1]
    pos = tf.reshape(tf.range(0.0, tf.to_float(T), dtype=tf.float32), [-1, 1])
    i = np.arange(0, repr_dim, 2, np.float32)
    denom = np.reshape(np.power(10000.0, i / repr_dim), [1, -1])
    enc = tf.expand_dims(tf.concat([tf.sin(pos / denom), tf.cos(pos / denom)],␣
 ↪1), 0)
    return tf.tile(enc, [tf.shape(inputs)[0], 1, 1]) * tf.expand_dims(tf.
 ↪to_float(mask), -1)

def label_smoothing(inputs, epsilon=0.1):
    C = inputs.get_shape().as_list()[-1]
    return ((1 - epsilon) * inputs) + (epsilon / C)

class Attention:
    def __init__(self, size_layer, embedded_size, learning_rate, size,␣
 ↪output_size,
                num_blocks = 2,
                num_heads = 8,
                min_freq = 50):
        self.X = tf.placeholder(tf.float32, (None, None, size))
        self.Y = tf.placeholder(tf.float32, (None, output_size))

        encoder_embedded = tf.layers.dense(self.X, embedded_size)
        encoder_embedded = tf.nn.dropout(encoder_embedded, keep_prob = 0.8)
        x_mean = tf.reduce_mean(self.X, axis = 2)
        en_masks = tf.sign(x_mean)
        encoder_embedded += sinusoidal_position_encoding(self.X, en_masks,␣
 ↪embedded_size)

        for i in range(num_blocks):
            with tf.variable_scope('encoder_self_attn_%d'%i,reuse=tf.
 ↪AUTO_REUSE):
                encoder_embedded = multihead_attn(queries = encoder_embedded,
                                        keys = encoder_embedded,
                                        q_masks = en_masks,
                                        k_masks = en_masks,
                                        future_binding = False,
```

```python
                                                num_units = size_layer,
                                                num_heads = num_heads)

            with tf.variable_scope('encoder_feedforward_%d'%i,reuse=tf.
 ↪AUTO_REUSE):
                encoder_embedded = pointwise_feedforward(encoder_embedded,
                                                embedded_size,
                                                activation = tf.nn.relu)

        self.logits = tf.layers.dense(encoder_embedded[-1], output_size)
        self.cost = tf.reduce_mean(tf.square(self.Y - self.logits))
        self.optimizer = tf.train.AdamOptimizer(learning_rate).minimize(
            self.cost
        )

def calculate_accuracy(real, predict):
    real = np.array(real) + 1
    predict = np.array(predict) + 1
    percentage = 1 - np.sqrt(np.mean(np.square((real - predict) / real)))
    return percentage * 100

def anchor(signal, weight):
    buffer = []
    last = signal[0]
    for i in signal:
        smoothed_val = last * weight + (1 - weight) * i
        buffer.append(smoothed_val)
        last = smoothed_val
    return buffer
```

```python
[7]: num_layers = 1
     size_layer = 128
     timestamp = 5
     epoch = 300
     dropout_rate = 0.8
     future_day = test_size
     learning_rate = 0.001
```

```python
[8]: def forecast():
         tf.reset_default_graph()
         modelnn = Attention(size_layer, size_layer, learning_rate, df_log.shape[1],␣
     ↪df_log.shape[1])
         sess = tf.InteractiveSession()
         sess.run(tf.global_variables_initializer())
         date_ori = pd.to_datetime(df.iloc[:, 0]).tolist()

         pbar = tqdm(range(epoch), desc = 'train loop')
```

```python
    for i in pbar:
        total_loss, total_acc = [], []
        for k in range(0, df_train.shape[0] - 1, timestamp):
            index = min(k + timestamp, df_train.shape[0] - 1)
            batch_x = np.expand_dims(
                df_train.iloc[k : index, :].values, axis = 0
            )
            batch_y = df_train.iloc[k + 1 : index + 1, :].values
            logits, _, loss = sess.run(
                [modelnn.logits, modelnn.optimizer, modelnn.cost],
                feed_dict = {
                    modelnn.X: batch_x,
                    modelnn.Y: batch_y
                },
            )
            total_loss.append(loss)
            total_acc.append(calculate_accuracy(batch_y[:, 0], logits[:, 0]))
        pbar.set_postfix(cost = np.mean(total_loss), acc = np.mean(total_acc))

    future_day = test_size

    output_predict = np.zeros((df_train.shape[0] + future_day, df_train.
↪shape[1]))
    output_predict[0] = df_train.iloc[0]
    upper_b = (df_train.shape[0] // timestamp) * timestamp

    for k in range(0, (df_train.shape[0] // timestamp) * timestamp, timestamp):
        out_logits = sess.run(
            modelnn.logits,
            feed_dict = {
                modelnn.X: np.expand_dims(
                    df_train.iloc[k : k + timestamp], axis = 0
                )
            },
        )
        output_predict[k + 1 : k + timestamp + 1] = out_logits

    if upper_b != df_train.shape[0]:
        out_logits = sess.run(
            modelnn.logits,
            feed_dict = {
                modelnn.X: np.expand_dims(df_train.iloc[upper_b:], axis = 0)
            },
        )
        output_predict[upper_b + 1 : df_train.shape[0] + 1] = out_logits
        future_day -= 1
        date_ori.append(date_ori[-1] + timedelta(days = 1))
```

6

```python
    for i in range(future_day):
        o = output_predict[-future_day - timestamp + i:-future_day + i]
        out_logits = sess.run(
            modelnn.logits,
            feed_dict = {
                modelnn.X: np.expand_dims(o, axis = 0)
            },
        )
        output_predict[-future_day + i] = out_logits[-1]
        date_ori.append(date_ori[-1] + timedelta(days = 1))

    output_predict = minmax.inverse_transform(output_predict)
    deep_future = anchor(output_predict[:, 0], 0.3)

    return deep_future[-test_size:]
```

```python
[9]: results = []
     for i in range(simulation_size):
         print('simulation %d'%(i + 1))
         results.append(forecast())
```

```
WARNING: Logging before flag parsing goes to stderr.
W0817 12:08:12.096583 140064997701440 deprecation.py:323] From <ipython-
input-6-24d2a24c36ef>:91: dense (from tensorflow.python.layers.core) is
deprecated and will be removed in a future version.
Instructions for updating:
Use keras.layers.dense instead.
W0817 12:08:12.104836 140064997701440 deprecation.py:506] From
/usr/local/lib/python3.6/dist-packages/tensorflow/python/ops/init_ops.py:1251:
calling VarianceScaling.__init__ (from tensorflow.python.ops.init_ops) with
dtype is deprecated and will be removed in a future version.
Instructions for updating:
Call initializer instance with the dtype argument instead of passing it to the
constructor

simulation 1

W0817 12:08:12.294501 140064997701440 deprecation.py:506] From <ipython-
input-6-24d2a24c36ef>:92: calling dropout (from tensorflow.python.ops.nn_ops)
with keep_prob is deprecated and will be removed in a future version.
Instructions for updating:
Please use `rate` instead of `keep_prob`. Rate should be set to `rate = 1 -
keep_prob`.
W0817 12:08:12.305350 140064997701440 deprecation.py:323] From <ipython-
input-6-24d2a24c36ef>:73: to_float (from tensorflow.python.ops.math_ops) is
deprecated and will be removed in a future version.
Instructions for updating:
```

Use `tf.cast` instead.
W0817 12:08:12.446460 140064997701440 deprecation.py:323] From <ipython-
input-6-24d2a24c36ef>:33: add_dispatch_support.<locals>.wrapper (from
tensorflow.python.ops.array_ops) is deprecated and will be removed in a future
version.
Instructions for updating:
Use tf.where in 2.0, which has the same broadcast rule as np.where
train loop: 100%|          | 300/300 [01:41<00:00,  2.97it/s, acc=96.7,
cost=0.00409]

simulation 2

train loop: 100%|          | 300/300 [01:40<00:00,  2.99it/s, acc=97.3,
cost=0.00184]

simulation 3

train loop: 100%|          | 300/300 [01:40<00:00,  2.98it/s, acc=96.7,
cost=0.00351]

simulation 4

train loop: 100%|          | 300/300 [01:40<00:00,  2.98it/s, acc=97.9,
cost=0.00112]

simulation 5

train loop: 100%|          | 300/300 [01:41<00:00,  2.97it/s, acc=98,
cost=0.00113]

simulation 6

train loop: 100%|          | 300/300 [01:40<00:00,  2.98it/s, acc=97.5,
cost=0.00165]

simulation 7

train loop: 100%|          | 300/300 [01:41<00:00,  2.96it/s, acc=95.8,
cost=0.00513]

simulation 9

train loop: 100%|          | 300/300 [01:41<00:00,  2.98it/s, acc=98,
cost=0.000974]

simulation 10

train loop: 100%|          | 300/300 [01:40<00:00,  2.99it/s, acc=96.8,
cost=0.00322]

```
[10]: accuracies = [calculate_accuracy(df['Close'].iloc[-test_size:].values, r) for r
      ↪in results]

      plt.figure(figsize = (15, 5))
      for no, r in enumerate(results):
```
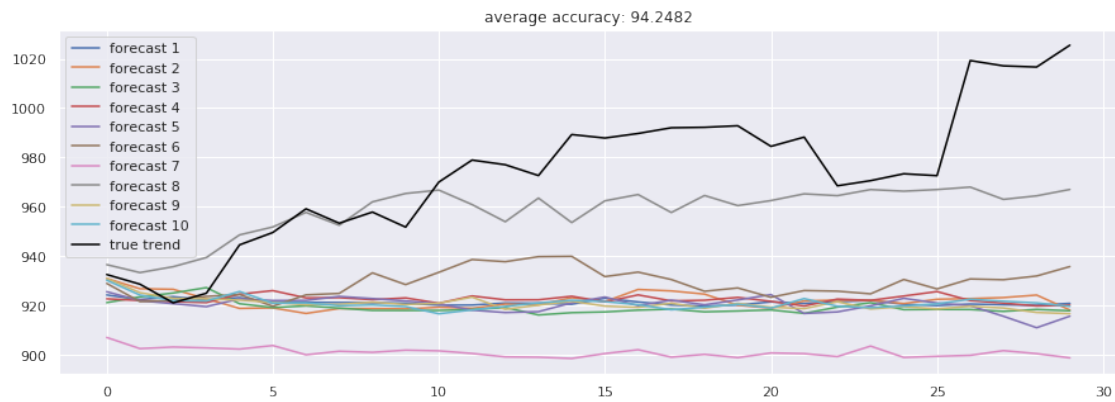
```
    plt.plot(r, label = 'forecast %d'%(no + 1))
plt.plot(df['Close'].iloc[-test_size:].values, label = 'true trend', c =␣
 ↪'black')
plt.legend()
plt.title('average accuracy: %.4f'%(np.mean(accuracies)))
plt.show()
```



[ ]: