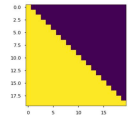


Wednesday, August 5, 2020 3:36 PM

- [The Annotated Transformer](#)
- [The Illustrated Transformer](#)
- [Attention Is All You Need](#)



```
def make_std_mask(tgt, pad):
    "Create a mask to hide padding and future words."
    tgt_mask = (tgt != pad).unsqueeze(-2)
    tgt_mask = tgt_mask & Variable(
        subsequent_mask(tgt.size(-1)).type_as(tgt_mask))
    return tgt_mask
```

```
def subsequent_mask(size):
    "Mask out subsequent positions."
    attn_shape = (1, size, size)
    subsequent_mask = np.triu(np.ones(attn_shape), k=1).astype(
        bool)
    return torch.from_numpy(subsequent_mask) == 0
```

```
def run_epoch(data_iter, model, loss_computer):
    "Standard Training and Logging Function"
    start = time.time()
    total_tokens = 0
    total_loss = 0
    tokens = 0
    for i, batch in enumerate(data_iter):
        out = model.forward(batch.src, batch.trg,
                             batch.src_mask, batch.trg_mask)
        loss, loss_computer(out, batch.trg_y, batch.ntokens)
        total_loss += loss
        total_tokens += batch.ntokens
        tokens += batch.ntokens
    if i % 50 == 1:
        elapsed = time.time() - start
        print("Epoch Step: %d Loss: %f Tokens per Sec: %f" %
              (i, total_loss / batch.ntokens, tokens / elapsed))
        start = time.time()
        tokens = 0
    return total_loss / total_tokens
```

```

class NeuronNet:
    """Neuron wrapper that implements rate"""
    def __init__(self, model_size, factor, warmup, optimizer):
        self.optimizer = optimizer
        self_step = 0
        self.warmup = warmup
        self.factor = factor
        self.model_size = model_size
        self.rate = 1

    def step(self):
        """update parameters and rate"""
        self_step += 1
        rate = self.rate
        for p in optimizer.param_groups:
            p["lr"] = rate
            self_step = rate
            self.optimizer.step()

    def train(self, step: None):
        """Implement 'rate' above"""
        if step is None:
            step = self_step
        return self.rate * \
            (self.model_size ** (-0.5)) * \
            min(step ** (-0.5), step * self.warmup ** (-1.5))

```

```

class LstmSmoothing(nn.Module):
    """Impement Lstm Smoothing"""
    def __init__(self, padding_idx, smoothing_p=0):
        super(LstmSmoothing, self).__init__()
        self.criterion = nn.KLDivLoss(size_average=False)
        self.padding_idx = padding_idx
        self.smoothing_confidence = 1.0 - smoothing_p
        self.smoothing = smoothing_p
        self.size = size
        self.true_dist = None

    def forward(self, x, target):
        assert x.size()[1] == self.size
        x.data = x.data.clone()
        true_dist, fill = self.smoothing (self.size - 2))
        true_dist.scatter_(1, target.data.unsqueeze(1), self.confidence)
        true_dist.fill_ (self.padding_idx, 1 - self.confidence)
        mask = torch.zeros(target.data == self.padding_idx)
        if mask.dim() > 0:
            true_dist.index_fill_(0, mask.squeeze(), 0.0)
        self.true_dist = true_dist
        return self.criterion(x, Variable(torch.ones(target.size(), 1)))

```

```
def greedy_decode(model, src, src_mask, max_len, start_symbol):
    memory, encoder_state = model.encode(src, src_mask)
    ys = torch.ones(1, 1).fill_(start_symbol).type_as(src.data)
    for i in range(max_len-1):
        out, decode_memory, src_mask, \
            variable(y), \
            variable(subsequent_mask(ys, size(1) \
                                     ).type_as(src.data)) = \
            model.generate(out, decode_memory, src_mask, \
                           variable(y), \
                           variable(subsequent_mask(ys, size(1) \
                                                       ).type_as(src.data)))
        prob = model.generator(out, -1)
        _, next_word = torch.max(prob, dim = 1)
        next_word = next_word.data[0]
        ys = torch.cat([ys, \
                        torch.ones(1, 1).type_as(src.data).fill_(next_word)], dim=1)
    return ys
```

```
model.eval()
src = Variable(torch.LongTensor([[1,2,3,4,5,6,7,8,9,10]]))
src_mask = Variable(torch.ones(1, 1, 10))
print(greedy_decode(model, src, src_mask, max_len=10, start_symbol=1))
```

```

def main(argv):
    """
    main function that will calculate total number of letters in a string"""
    if len(argv) != 2:
        print("Usage: %s <string>" % argv[0])
        return

    str = argv[1]

    # calculate the number of letters in the string
    num_letters = 0
    for char in str:
        if char.isalpha():
            num_letters += 1

    print("Number of letters in '%s' is: %d" % (str, num_letters))

if __name__ == "__main__":
    main(sys.argv)

```

Code Snippet 2: A Python script to calculate the number of letters in a string.

```

def main(argv):
    """
    main function that will calculate total number of letters in a string"""
    if len(argv) != 2:
        print("Usage: %s <string>" % argv[0])
        return

    str = argv[1]

    # calculate the number of letters in the string
    num_letters = 0
    for char in str:
        if char.isalpha():
            num_letters += 1

    print("Number of letters in '%s' is: %d" % (str, num_letters))

if __name__ == "__main__":
    main(sys.argv)

```

Code Snippet 3: A Python script to calculate the number of letters in a string, using a more concise approach.

```

def main(argv):
    """
    main function that will calculate total number of letters in a string"""
    if len(argv) != 2:
        print("Usage: %s <string>" % argv[0])
        return

    str = argv[1]

    # calculate the number of letters in the string
    num_letters = sum(1 for char in str if char.isalpha())

    print("Number of letters in '%s' is: %d" % (str, num_letters))

if __name__ == "__main__":
    main(sys.argv)

```

Code Snippet 4: A Python script to calculate the number of letters in a string, using a more concise approach, with a more descriptive variable name.

```

def main(argv):
    """
    main function that will calculate total number of letters in a string"""
    if len(argv) != 2:
        print("Usage: %s <string>" % argv[0])
        return

    str = argv[1]

    # calculate the number of letters in the string
    num_letters = sum(1 for char in str if char.isalpha())

    print("Number of letters in '%s' is: %d" % (str, num_letters))

if __name__ == "__main__":
    main(sys.argv)

```

Code Snippet 5: A Python script to calculate the number of letters in a string, using a more concise approach, with a more descriptive variable name, and a more descriptive docstring.

```

def main(argv):
    """
    main function that will calculate total number of letters in a string"""
    if len(argv) != 2:
        print("Usage: %s <string>" % argv[0])
        return

    str = argv[1]

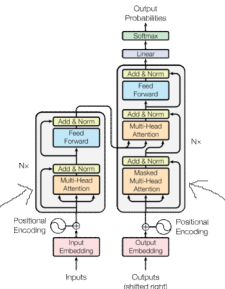
    # calculate the number of letters in the string
    num_letters = sum(1 for char in str if char.isalpha())

    print("Number of letters in '%s' is: %d" % (str, num_letters))

if __name__ == "__main__":
    main(sys.argv)

```

```
def clones(module, N):
    "Produce N identical layers."
    return nn.ModuleList([copy.deepcopy(module) for _ in range(N)])
```

[illegible]

```
class EncoderLayer(nn.Module):
    "Encoder is made up of self-attn and
    def __init__(self, size, self_attn, f
        super(EncoderLayer, self).__init
        self.self_attn = self_attn
```

```
self.feed_dict = feed_dict
self.feed_forward = feed_forward
self.sublayer = clones(SublayerC)
self.size = size
```

```
def forward(self, x, mask):
    "Follow Figure 1 (left) for connection"
    x = self.sublayer[0](x, lambda x:
    return self.sublayer[1](x, self.
```

```
class DecoderLayer(nn.Module):
    "Decoder is made of self-attn, src-attn"
```

```
def __init__(self, size, self_attn, src
super(DecoderLayer, self).__init__(
self.size = size
self.self_attn = self_attn
self.src_attn = src_attn
```

```

self.feed_forward = feed_forward
self.sublayer = clones(SublayerComm

def forward(self, x, memory, src_mask,
    "Follow Figure 1 (right) for connec

```

```
m = memory
x = self.sublayer[0](x, lambda x: s
x = self.sublayer[1](x, lambda x: s
return self.sublayer[2](x, self.fee
```

```
class LayerNorm(nn.Module):
    "Construct a layernorm module (See citation for details)."
    def __init__(self, features, eps=1e-6):
        super(LayerNorm, self).__init__()
        self.a_2 = nn.Parameter(Torch.ones(features))
        self.b_2 = nn.Parameter(Torch.zeros(features))
        self.eps = eps

    def forward(self, x):
        mean = x.mean(-1, keepdim=True)
        std = x.std(-1, keepdim=True)
        return self.a_2 * (x - mean) / (std + self.eps) + self.b_2
```

```
class EncoderLayer(nn.Module):
    "Encoder is made up of self-attn and feed forward (defi
    def __init__(self, size, self_attn, feed_forward, dropo
        super(EncoderLayer, self).__init__()
        self.self_attn = self_attn
        self.feed_forward = feed_forward
        self.sublayer = clones(SublayerConnection(size, dro
        self.size = size
```

```
def forward(self, x, mask):
    """Follow Figure 1 (left) for connections."""
    x = self.sublayer[0](x, lambda x: self.self_attn(x,
        return self.sublayer[1](x, self.feed_forward)

class DecoderLayer(nn.Module):
    """Decoder is made of self-attn, src-attn, and feed forward"""
```

```
def __init__(self, size, self_attn, src_attn, feed_forward,
             super(DecoderLayer, self).__init__())
    self.size = size
    self.self_attn = self_attn
    self.src_attn = src_attn
    self.feed_forward = feed_forward
    self.sublayer = clones(SublayerConnection(size, dropout),
                             2)

def forward(self, x, memory, src_mask, tgt_mask):
    "Follow figure 1 (right) for connections."
```

```
m = memory
x = self.sublayer[0](x, lambda x: self.self_attn(x, x,
x = self.sublayer[1](x, lambda x: self.src_attn(x, m,
return self.sublayer[2](x, self.feed_forward)
```

4. To counteract this effect, we scale the dot products by $\frac{1}{\sqrt{d_k}}$.

```
Attention(Q, K, V) = softmax( $\frac{QK^T}{\sqrt{d_k}}$ )V
```

```
def attention(query, key, value, mask=None, dropout=None):
    "Compute 'Scaled Product Attention' on three inputs"
    d_k = query.size(-1)
    attn = torch.matmul(query, key.transpose(-2, -1)) \
           / math.sqrt(d_k)
    if mask is not None:
        scores = scores.masked_fill(mask == 0, -1e9)
    p_attn = F.softmax(scores, dim=-1)
    if dropout is not None:
        p_attn = dropout(p_attn)
    return torch.matmul(p_attn, value), p_attn
```

```

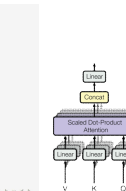
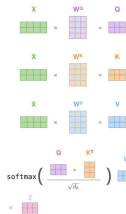
class LinearTransformation(nn.Module):
    def __init__(self, n, d_model, dropout):
        """ In model size and number of heads. """
        super(LinearTransformation).__init__()
        assert d_model % n == 0
        # we assume d_model always equals d_B
        self.d_B = d_model // n
        self.n = n
        self.linear = clones(nn.Linear(n*d_B, d_model), 4)
        self.dropout = nn.Dropout(dropout)

    def forward(self, query, key, value, mask=None):
        """ Implements Trm in:
        # If not in test mode:
        #   Same mask applied to all 4 heads.
        #   mask = mask.unsqueeze(-2)
        #   n_batches = query.size(0)
        #   return self.linear(query)"""
        # self all the linear projections in batch from n * B to B * n * d_B
        query, key, value = [
            v.view(-1, self.n, self.d_B, transpose(2))
            for v in [query, key, value]]
        # self attn operates on all the projected tensors in batch
        query, attn = attention(query, key, value, mask,
                                dropout=self.dropout)
        # self "Concat" using a view and apply a final linear.
        x = transpose(1, 2)concaten([
            v.view(-1, self.n, self.d_B)
            for v in [query, attn]])
        return self.linear(x)

```

```
class SublayerConnection(nn.Module):
    """
    A residual connection followed by a layer norm.
    Note for code simplicity the norm is first as opposed to last.
    """
    def __init__(self, size, dropout):
        super(SublayerConnection, self).__init__()
        self.norm = LayerNorm(size)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x, sublayer):
        """
        Apply residual connection to any sublayer with the same size."
        """
        return x + self.dropout(sublayer(self.norm(x)))
```



```
class Generator(nn.Module):
    "Define standard linear + softmax generation step."
    def __init__(self, d_model, vocab):
        super(Generator, self).__init__()
        self.proj = nn.Linear(d_model, vocab)

    def forward(self, x):
        return F.log_softmax(self.proj(x), dim=-1)
```

```
class PositionalEncoding(nn.Module):
    "Implement the PE function."
    def __init__(self, d_model, dropout, max_len=5000):
        super(PositionalEncoding, self).__init__()
        self.dropout = nn.Dropout(dropout)

        # Compute the positional encodings once in log space.
        pe = torch.zeros(max_len, d_model)
        position = torch.arange(0, max_len).unsqueeze(1)
        div_term = torch.exp(torch.arange(0, d_model, 2) *
                              -(math.log(10000.0) / d_model))
        pe[:, 0::2] = torch.sin(position * div_term)
        pe[:, 1::2] = torch.cos(position * div_term)
        pe = pe.unsqueeze(0)
        self.register_buffer('pe', pe)

    def forward(self, x):
        x = x + Variable(self.pe[:, :x.size(1)]),
            requires_grad=False)
        return self.dropout(x)
```



```
m = memory
x = self.sublayer[0](x, lambda x: self.self_attn(x, x, x, tgt_mask))
x = self.sublayer[1](x, lambda x: self.src_attn(x, m, m, src_mask))
return self.sublayer[2](x, self.feed_forward)
```

```
class PositionwiseFeedForward(nn.Module):
    "Implements FFN equation."
    def __init__(self, d_model, d_ff, dropout=0.1):
        super(PositionwiseFeedForward, self).__init__()
        self.w_1 = nn.Linear(d_model, d_ff)
        self.w_2 = nn.Linear(d_ff, d_model)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x):
        return self.w_2(self.dropout(F.relu(self.w_1(x))))
```

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$